

HIGHER ORDER TIME FILTERS FOR EVOLUTION EQUATIONS

by

Ahmet Guzel

Bachelor of Science in Mathematics,
Uludag University, Bursa, Turkey, 2008

Master of Science in Mathematics,
University of Texas at San Antonio, Texas, 2012

Master of Arts in Mathematics,
University of Pittsburgh, Pittsburgh, Pennsylvania, 2016

Submitted to the Graduate Faculty of
the Kenneth P. Dietrich School of Arts and Sciences
in partial fulfillment

of the requirements for the degree of

Doctor of Philosophy

University of Pittsburgh

2018

UNIVERSITY OF PITTSBURGH
KENNETH P. DIETRICH SCHOOL OF ARTS AND SCIENCES

This dissertation was presented

by

Ahmet Guzel

It was defended on

March 20, 2018

and approved by

Dr. Catalin Trenchea, Associate Professor, Dept. of Mathematics

Dr. William Layton, Professor, Dept. of Mathematics

Dr. Michael Neilan, Associate Professor, Dept. of Mathematics

Dr. Patrick Smolinski, Associate Professor, Dept. of Mechanical Engineering

Dissertation Director: Dr. Catalin Trenchea, Associate Professor, Dept. of Mathematics

Copyright © by Ahmet Guzel
2018

HIGHER ORDER TIME FILTERS FOR EVOLUTION EQUATIONS

Ahmet Guzel, PhD

University of Pittsburgh, 2018

Time filter is a non-intrusive technique that post-processes the previously computed values of given numerical methods. The purpose of this study is to construct new time filters, which will increase the accuracy and stability of existing legacy codes. We focus on time filters for the leapfrog method and the backward Euler method.

The leapfrog scheme is a second-order, symplectic, explicit method, which is widely used in the numerical models of weather and climate, currently in conjunction with the Robert-Asselin (RA) and Robert-Asselin-Williams (RAW) time filters.

- We propose and analyze a novel filter, which combines the higher-order Robert-Asselin (hoRA) filter with a Williams' step (LF-hoRAW). This filter better addresses the issue of time-splitting instability of the leapfrog scheme and increases the stability of hoRA, reduces the magnitude of the truncation error, improves the accuracy of amplitude compared to the hoRA, and conserves the three-time-level mean.
- We perform linear error analysis for general high-order Robert-Asselin (ghoRA) time filter applied to the leapfrog scheme, and derive the phase and amplitude errors for a pre-determined order of accuracy using the modified equation.

The fully implicit (backward) Euler method is one of the first method commonly implemented when extending a code for the steady state problem, and often the method of last resort for complex applications.

- We construct a time filter for the backward Euler method, which reduces the discrete curvature of the solution, increases accuracy from first to second-order, gives an immediate error estimator and induces an equivalent two-step, A -stable, linear multistep method.

Keywords: A -stability, Backward Euler method, Leapfrog method, Modified equation, Robert-Asselin-Williams, Time filters.

TABLE OF CONTENTS

PREFACE	xii
1.0 INTRODUCTION	1
2.0 HIGHER ORDER ROBERT-ASSELIN-WILLIAMS TIME FILTER	5
2.1 Linear Analysis	6
2.1.1 Previous Work	7
2.1.2 The LF-hoRAW as a Linear Multistep Method	8
2.1.3 The Consistency Order of LF-hoRAW	11
2.1.4 The Stability Domain of LF-hoRAW	11
2.2 Curvature Evolution	13
2.3 Error Analysis for Phase and Amplitude	15
2.4 Comparison of LF-hoRA, LF-hoRAW and AB3 Methods	19
2.5 Numerical Tests	22
2.5.1 Simple Pendulum	22
2.5.2 Ozone Photochemistry	24
2.5.3 Lorenz System	25
2.6 Summary	27
3.0 THE GENERAL HIGH-ORDER ROBERT-ASSELIN TIME FILTER	28
3.1 Previous Work	28
3.2 Error Analysis	30
3.3 Summary	35
4.0 BACKWARD EULER PLUS FILTER	36
4.1 Constant Time Step	38

4.1.1	Derivation of the Method	39
4.1.2	Stability for Constant Time Step	40
4.2	Variable Time Step	42
4.2.1	Time Filter for Variable Time Step	42
4.2.2	The Local Truncation Error	44
4.2.3	Stability for Variable Time Step	46
4.2.4	Adaptive Time Step Algorithm	48
4.3	Error Analysis for Phase and Amplitude	49
4.4	Numerical Tests	53
4.4.1	The Lorenz System	53
4.4.2	Preservation of Lyapunov Stability	54
4.4.3	Periodic and Quasi-Periodic Oscillations	55
4.4.4	The Van der Pol Equation	57
4.5	Summary	59
5.0	CONCLUSIONS	60
	BIBLIOGRAPHY	61

LIST OF TABLES

2.1	Summary of the conservation of three-time-level mean, stability, and accuracy properties of the LF-hoRAW for some values of α	19
2.2	The comparison of LF-hoRAW, LF-hoRA and AB3 schemes with some featured values of α and β	22
4.1	The comparison of halving, doubling and the same step using variable time step backward Euler and backward Euler plus filter algorithm for the Van der Pol equation.	58

LIST OF FIGURES

2.1	The exact solution of simple harmonic motion for variable $x(t)$ with four numerical solutions using $\Delta t = 0.2$ s	7
2.2	The amplification factors of the physical mode (solid line) and two computational modes (dotted line) of LF-hoRAW	10
2.3	The magnitudes of the physical mode (solid line) and computational modes (dotted line) of LF-hoRAW.	10
2.4	Root locus curve of LF-hoRAW with various α and β	12
2.5	The magnitude of physical mode amplitudes while $\alpha \gtrapprox \frac{2-\beta}{8-5\beta}$ for $\beta = 0.2$ (left) and $\beta = 0.4$ (right).	13
2.6	The hoRAW filter moves the inner and right outer points through displacements $\alpha\beta(\kappa_{\text{old}}^n - \kappa_{\text{new}}^{n-1})/2$ and $(1 - \alpha)\beta(\kappa_{\text{old}}^n - \kappa_{\text{new}}^{n-1})/2$, respectively	14
2.7	Amplitude (top) and relative phase change (bottom) of the physical mode of LF-hoRAW	18
2.8	Root locus curves of LF-hoRA, LF-hoRAW the AB3 schemes	20
2.9	Comparison of the coefficients $C_2^{1\beta}, C_2^{\alpha\beta}$ (2.8) and $C_{AB3} = 3/8$	20
2.10	The comparison of amplitude (top) and relative phase (bottom) of physical mode of LF-hoRA, LF-hoRAW and AB3.	21
2.11	Numerical solution to the simple pendulum problem	23
2.12	Numerical solutions for chemical concentrations	25
2.13	Computed numerical solutions to the Lorenz system	26
4.1	Stability region of backward Euler plus time filter	41
4.2	Boundaries of Stability Regions	42

4.3	A -stable for $\nu \leq$ dark curve, Dashed Curve = $\mathcal{O}(\Delta t^2)$	47
4.4	Comparison of phase speed and amplitude of physical mode for backward Euler method and second-order backward Euler plus filter method.	53
4.5	Numerical solution of X for the Lorenz system with time step $\Delta t = 0.01$ (left) and $\Delta t = 0.02$ (right).	54
4.6	Numerical solution of Lyapunov stability test problem with time step $\Delta t = 0.2$	55
4.7	Numerical solution of pendulum test problem with time step $\Delta t = 0.1$	56
4.8	Exact soln, non-adaptive and adaptive backward Euler plus filter	57
4.9	Numerical solution of Van der Pol equation using variable time step backward Euler and backward Euler plus filter	58

LIST OF ALGORITHMS

4.1 Halving and Doubling Time step 48

PREFACE

I would like to express my deepest gratitude and appreciation to my advisor, Professor Catalin Trenchea, for his unwavering support, collegiality and mentorship throughout my graduate study. I would also like to thank him for always taking his time and providing answers to my endless questions towards my research. Besides mathematics, He is a great friend and life mentor. It is a great honor for me to have him as my advisor.

I would like to extend my gratitude to Professor William Layton for his continuous support, encouragement and priceless guidance on my research. His deep knowledge and a unique perspective on research and mathematics helped me to make this dissertation better.

I would also like to thank Professor Michael Neilan and Professor Patrick Smolinski for giving their valuable time to be part of my committee. I am also grateful for their useful comments and suggestions.

I wish to thank all my fellow numerical analysis group, staff and faculty in Department of Mathematics for the help and support they provided during my graduate study.

I would also like to express my sincerest gratitude and appreciation to my parents, Halil and Menci, my wife, Fatima Zehra, my siblings, and their families for moral support, endless patience and constant encouragement during my studies in the United States.

I would like to thank my daughters, Firdevs and Cennet, for always making me smile and understanding when I was at school instead of spending time with them.

I dedicate this dissertation to my daughters, Firdevs and Cennet, with love.

1.0 INTRODUCTION

The computational approach for predicting the behaviour of a dynamical system consists in the construction of an algorithm that accurately simulates its evolution. The task of the forecasting procedure may be reduced to the following iterative algorithm: given the current state of the system (the input), use the governing equations to approximate the state at a slightly later time (the output), and repeat as many times as needed. The output of this algorithm exhibits *modeling errors* (unknown initial and boundary data, forcing terms, and the numerical model) and *numerical approximation errors* (parameterizations of sub-grid processes, semi-discretization in time and space). In this thesis we focus on reducing the errors due to the semi-discrete time approximation. Specifically, we develop and analyze post-processing non-intrusive techniques to obtain more accurate solutions, at low computational costs, from solutions generated by existing numerical algorithms.

Consider the initial value problem (IVP)

$$u'(t) = f(u(t)), \text{ for } t > 0 \text{ and } u(0) = u_0 \quad (1.1)$$

which is approximated by the following k -step ($k \geq 1$) method

$$\sum_{i=0}^k \alpha_i u_{n+1-i} = \Delta t f\left(\sum_{i=0}^k \beta_i u_{n+1-i}\right), \text{ given } \{u_i\}_{i=0}^{k-1}. \quad (1.2)$$

The relation (1.2) defines a one-leg linear multistep method [12], where $\alpha_0 \neq 0$, $\sum_{i=0}^k \beta_i = 1$ and the starting values $\{u_i\}_{i=0}^{k-1}$ are given or computed with a different numerical method¹. Here Δt denotes the constant time step, u_n is the numerical solution approximating the

¹The linear multistep method (1.2) is explicit if $\beta_0 = 0$, otherwise it is implicit.

exact solution $u(t_n)$ at time $t_n = n \Delta t$. This notation convention will be used through out the text. We define a time filter for the linear multistep method (1.2) as follows:

$$\overline{u_{n+\theta}} = \gamma_0 u_{n+1} + \gamma_1 u_n + \sum_{j=2}^{\ell} \gamma_j \overline{u_{n+1-j}}, \quad \ell \geq k \quad (1.3)$$

where $\overline{u_{n+\theta}}$ denotes the once filtered solution, approximating $u_{n+\theta}$, using previously computed values, γ_j are real numbers and θ is either 0 or 1.

In this work we focus on time filters for two widely used numerical methods applied to the first-order differential equation, namely, the leapfrog and backward Euler method.

- The leapfrog scheme applied to (1.1) is given by

$$u_{n+1} - u_{n-1} = 2\Delta t f(u_n), \text{ given } u_0, u_1. \quad (1.4)$$

In Chapter 2, we introduce and analyze the leapfrog scheme filtered twice with a higher-order Robert-Asselin time filter and Williams' step. In Chapter 3, the leapfrog scheme is filtered once by the general high-order Robert-Asselin time filter (LF-ghoRA) given as

$$\overline{u_n} = u_n + a_{n+1}u_{n+1} + a_n u_n + \sum_{j=1}^k a_{n-j} \overline{u_{n-j}}, \quad k \geq 1.$$

- The backward Euler method applied to (1.1) is given by

$$u_{n+1} - u_n = \Delta t f(u_{n+1}), \text{ given } u_0. \quad (1.5)$$

In Chapter 4, the solution of backward Euler method is filtered once by a simple finite difference approximation of the second time-derivative given as

$$\overline{u_{n+1}} = u_{n+1} - \frac{\nu}{2}(u_{n+1} - 2\overline{u_n} + \overline{u_{n-1}}).$$

The filtered solution \overline{u} of the combination of the given numerical method (1.2) with the time filter (1.3) exhibits better stability and accuracy compared to unfiltered solution u . In the following, we give a brief account of existing time filters for the leapfrog scheme, currently used in the weather and climate models.

The leapfrog scheme, also known as the midpoint rule or the explicit Nyström method, is an explicit, two-step, three-time-level, second-order accurate, weakly-stable, neutral time

stepping scheme. It is best suited for the time integration of linear oscillatory systems and is widely used in weather and climate computational models. The major weakness of the leapfrog scheme is the spurious growth of the computational mode when applied to nonlinear equations [16, 34, 49], the so-called “time-splitting” instability [15, 33, 48]. There are various ways to damp computational mode of leapfrog scheme, see e.g. [16]. In the atmospheric science, it is common to control the computational mode by non-intrusively post-process the leapfrog scheme based legacy codes through a second-order time filter. This filter is closely related to the centered second-derivative time filter

$$\overline{u}_n = u_n + \gamma(u_{n+1} - 2u_n - u_{n-1}),$$

where u_n denotes the solution at time $n\Delta t$ prior to time filtering, \overline{u}_n is the solution after filtering and γ is a positive real constant which determines the strength of the filter.

The Robert-Asselin (RA) time filter, designed by Robert [30] and analyzed by Asselin [2], filters once the middle value u_n obtained by (1.4) into

$$\overline{u}_n = u_n + \frac{\nu}{2}(u_{n+1} - 2u_n + \overline{u}_{n-1}).$$

The combination of LF-RA successfully suppresses the leapfrog scheme’s computational mode, but also weakly damps the physical mode, reducing the second-order accuracy of the unfiltered leapfrog scheme to first-order. The effect of RA filter has been investigated in [4, 7, 14, 41]. It is currently used in the majority of the operational numerical weather prediction models, atmospheric general circulation models for climate simulation, ocean general circulation models, models of the fluids in rotating annulus laboratory experiments (see e.g., [48] and references therein).

Williams [1, 48, 49] made a significant improvement to the RA filter, altering both values u_n and u_{n+1} , obtained by (1.4) into

$$\begin{aligned}\overline{\overline{u}}_n &= \overline{u}_n + \frac{\nu\alpha}{2}(u_{n+1} - 2\overline{u}_n + \overline{\overline{u}}_{n-1}), \text{ and} \\ \overline{\overline{u}}_{n+1} &= u_{n+1} - \frac{\nu(1-\alpha)}{2}(u_{n+1} - 2\overline{u}_n + \overline{\overline{u}}_{n-1})\end{aligned}$$

respectively. Filtering one more time compared to RA, the RAW filtered leapfrog scheme almost conserves the three-time-level mean of the predicted field, increases the accuracy of

amplitude errors by two orders, yielding third-order accuracy, and greatly reduces the magnitude of the first-order truncation error. The RAW filter has been studied and implemented in various model (see [38, 40, 47, 51, 52]).

Using a filter closely related to the third time-derivative, Li and Trenchea [33] introduced a higher-order Robert-Asselin (hoRA) type time filter. hoRA filters once the middle value u_n obtained by (1.4) into

$$\overline{u_n} = u_n + \frac{\beta}{2}(u_{n+1} - 2u_n + \overline{u_{n-1}}) - \frac{\beta}{2}(u_n - 2\overline{u_{n-1}} + \overline{u_{n-2}}).$$

It is a linear post-process to the leapfrog scheme, which controls the undamped computational modes, and increases the numerical accuracy of the RA time filter to third-order when $\beta = 0.4$, yielding fourth-order accuracy for the phase and amplitude of the physical mode [33, 34].

The plan of this thesis as follows. In Chapter 2, we propose an extension of the hoRA time filter, by altering both values u_n, u_{n+1} obtained by (LF) with a hoRA step and a Williams-type step. This combination further increases the stability, reduces the magnitude of the truncation error, improves the accuracy of amplitude compared to the hoRA filtered leapfrog scheme, and also conserves the three-time-level mean. In Chapter 3 we perform the error analysis of the general high-order Robert-Asselin (ghoRA) time filter using a modified equations. The modified equation gives a natural and simple way to find the error distribution between phase and amplitude. In Chapter 4 we construct a time filter for the backward Euler method, and show that the combination reduces the discrete curvature of the solution, increases the accuracy from first to second-order, gives an immediate error estimator and induces a method akin to BDF2. The effect of each step in the combination of 2-step is conceptually clear. The combination also extends easily to variable time steps.

2.0 HIGHER ORDER ROBERT-ASSELIN-WILLIAMS TIME FILTER

In this chapter, we construct and analyzed a higher order Robert-Asselin-Williams(hoRAW) time filter. The work of this chapter is based on [23]. The proposed hoRAW filtered leapfrog scheme applied to initial value problem (1.1) is:

$$\begin{aligned} w_{n+1} &= u_{n-1} + 2\Delta t f(v_n) \\ u_n &= v_n + \frac{\alpha\beta}{2}(w_{n+1} - 2v_n + u_{n-1}) - \frac{\alpha\beta}{2}(v_n - 2u_{n-1} + u_{n-2}) \\ v_{n+1} &= w_{n+1} + \frac{\beta(\alpha-1)}{2}(w_{n+1} - 2v_n + u_{n-1}) - \frac{\beta(\alpha-1)}{2}(v_n - 2u_{n-1} + u_{n-2}) \end{aligned}$$

where the dimensionless parameters $\beta \in [0, 1]$ and $\alpha \in [0, 1]$. Here w , v , u denote the unfiltered, once and twice filtered values, respectively. The last two terms in each step can be combined as $w_{n+1} - 3v_n + 3u_{n-1} - u_{n-2}$, which is a finite difference approximation to the third time-derivative. The LF-hoRAW is generally second-order accurate, and third-order when $\alpha = \frac{2+2\beta}{7\beta}$, yielding fourth-order accuracy for the phase and amplitude of the physical mode. Using a backward error analysis approach, the modified equation, we shall prove that when $\alpha = \frac{2-\beta}{8-5\beta}$, the LF-hoRAW method achieves sixth-order accuracy in amplitude. The hoRAW time filter has a twenty percent increase in stability compared hoRA, and the LF-hoRAW is twenty five percent more stable than the AB3 method:

$$u_{n+1} = u_n + \frac{\Delta t}{12}(23f(u_n) - 16f(u_{n-1}) + 5f(u_{n-2}))$$

when

$$\alpha = \frac{4 - 12\beta + 5\beta^2 - 2\sqrt{4 + 12\beta - 15\beta^2 + 4\beta^3}}{25\beta^2 - 36\beta}.$$

The storage factor for the leapfrog scheme combined with hoRAW filter is five. Compared with the intrusive AB3 method (three function evaluations per time iteration), the hoRAW

filtered leapfrog scheme is almost as accurate, stable and efficient, yet non-intrusive and easily implementable in existing legacy codes. We briefly illustrate the improvement that may be achieved by the proposed hoRAW filter, when used in conjunction with the leapfrog scheme, by numerically integrating the simple harmonic motion

$$\begin{aligned}\frac{dx}{dt} &= -y, & x(0) &= 1, \\ \frac{dy}{dt} &= x, & y(0) &= 0.\end{aligned}$$

We compare the exact solution with the numerical solutions of LF-RA ($\nu = 0.2$), LF-RAW ($\nu = 0.2, \alpha = 0.53$), LF-hoRA ($\beta = 0.1$) and LF-hoRAW ($\beta = 0.1, \alpha = 0.27$)¹ in Figure 2.1. The amplitude error of the hoRAW filter is significantly smaller than the RA, RAW and hoRA filters². The energy of the oscillation corresponding to $x^2 + y^2$, which is conserved by the continuous equations, decreases to 0 using the RA filter, decreases to 57% using the RAW filter, to 70% using hoRA, but is 99% conserved using the hoRAW filter, on the time interval $[0, 500]$.

2.1 LINEAR ANALYSIS

The phase and amplitude errors of time stepping schemes for non-dissipative dynamical systems is typically evaluated by analyzing the solutions of the oscillation equation (see [16, 24]),

$$u'(t) = i\omega u(t) \tag{2.1}$$

where ω is real constant. In this section we derive the consistency and stability properties of the hoRAW filter. First we briefly recall the properties of the Robert-Asselin, Robert-Asselin-Williams and hoRA time filters.

¹The chosen values yield, in the limit of good time resolution, the same damping rates of computational modes (The damping rate of the computational mode for both LF-RA and LF-RAW is $1 - \nu$. The damping rates of the most unstable computational mode of LF-hoRA is $1 - 2\beta$. The damping rate of the most unstable computational mode of LF-hoRAW is $\frac{2-3\beta-\alpha\beta+\sqrt{4+4\beta(\alpha-5)+(\alpha\beta+3\beta)^2}}{4}$).

²We note that for these parameter values, LF-RA has second-order amplitude accuracy, LF-RAW is almost fourth-order, LF-hoRA is fourth-order, while LF-hoRAW is sixth-order accurate in amplitude.

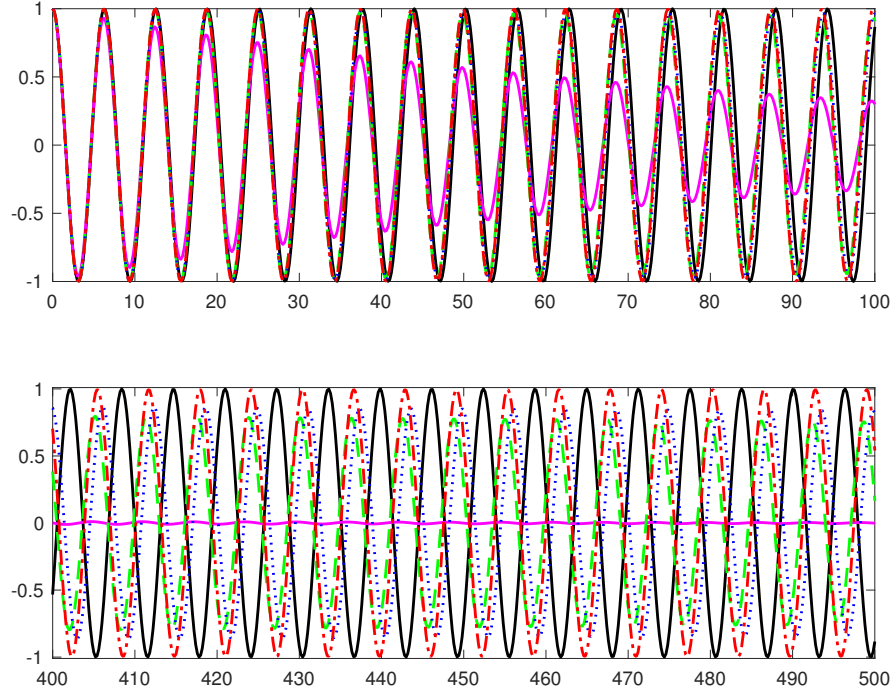


Figure 2.1: The exact solution of simple harmonic motion for variable $x(t)$ with four numerical solutions using $\Delta t = 0.2$ s, on the time interval $[0, 100]$, and $[400, 500]$.

(Exact —, LF-RA —, LF-hoRA ···, LF-RAW ---, and LF-hoRAW -.-)

2.1.1 Previous Work

The RAW filtered leapfrog scheme applied to (2.1) writes

$$w_{n+1} = u_{n-1} + 2i\omega\Delta t v_n, \quad (\text{Leapfrog})$$

$$u_n = v_n + \frac{\alpha\nu}{2}(w_{n+1} - 2v_n + u_{n-1}), \quad (\text{Robert-Asselin})$$

$$v_{n+1} = w_{n+1} + \frac{(\alpha - 1)\nu}{2}(w_{n+1} - 2v_n + u_{n-1}), \quad (\text{Williams})$$

where w , v , u are the unfiltered, once filtered and twice filtered values, respectively. The dimensionless parameters $\nu \in [0, 1]$ and $\alpha \in [0.5, 1]$. When $\alpha = 1$ the (Williams) step drops

out and the LF-RAW becomes the LF-RA scheme, and when $\nu = 0$ the leapfrog scheme is recovered. Both RA and RAW filters are generally first-order accurate and successfully dampen the computational mode. However, the RAW filter provides a higher accuracy for the amplitude of the physical mode, compared to the RA filtered leapfrog scheme. When $\alpha = 0.5$, the LF-RAW preserves three-time-level mean, it is second-order accurate, yielding third-order accuracy for the amplitude of the physical mode; however LF-RAW is unconditionally unstable in this case. Nevertheless, with α slightly larger than 0.5, e.g., $\alpha = 0.53$, LF-RAW yields almost third-order accuracy for the amplitude of the physical mode (see [34, 48]).

The hoRA filtered leapfrog (LF-hoRA) applied to (2.1) is given by

$$\begin{aligned} v_{n+1} &= u_{n-1} + 2i\omega\Delta t v_n & (\text{Leapfrog}) \\ u_n &= v_n + \frac{\beta}{2}(v_{n+1} - 2v_n + u_{n-1}) - \frac{\beta}{2}(v_n - 2u_{n-1} + u_{n-2}) & (\text{high-order RA}) \end{aligned}$$

where v , u are the unfiltered and once filtered values, respectively, and $\beta \in (0, 1)$. In the limit of good time resolution, i.e., $\omega\Delta t \ll 1$, the LF-hoRA scheme is generally second-order accurate, and third-order accurate when $\beta = 0.4$ (see [33]).

2.1.2 The LF-hoRAW as a Linear Multistep Method

The hoRAW filtered leapfrog(LF-hoRAW) scheme applied (2.1) writes as the following

$$\begin{aligned} w_{n+1} &= u_{n-1} + 2i\omega\Delta t v_n & (\text{LF}) \\ u_n &= v_n + \frac{\alpha\beta}{2}(w_{n+1} - 2v_n + u_{n-1}) - \frac{\alpha\beta}{2}(v_n - 2u_{n-1} + u_{n-2}) & (\text{hoRA}) \\ v_{n+1} &= w_{n+1} + \frac{\beta(\alpha-1)}{2}(w_{n+1} - 2v_n + u_{n-1}) - \frac{\beta(\alpha-1)}{2}(v_n - 2u_{n-1} + u_{n-2}) & (\text{W}) \end{aligned}$$

where w , v , u are unfiltered, once filtered and twice filtered values, respectively and dimensionless parameter $\beta \in [0, 1]$ and $\alpha \in [0, 1]$ ³.

³The LF-hoRA is recovered when $\alpha = 1$.

First we solve the linear system (LF)-(hoRA)-(W) for w_{n+1} , v_n , v_{n+1} in terms of u_n, u_{n-1}, u_{n-2} , we obtain

$$v_n = \frac{u_n - 2\alpha\beta u_{n-1} + \frac{\alpha\beta}{2}u_{n-2}}{1 - \frac{3\alpha\beta}{2} + i\omega\Delta t\alpha\beta},$$

$$v_{n+1} = (1 + 2\alpha\beta - 2\beta)u_{n-1} - \frac{\alpha\beta - \beta}{2}u_{n-2}$$

$$+ \frac{4i\omega\Delta t + 2i\omega\Delta t\alpha\beta - 2i\omega\Delta t\beta - 3\alpha\beta + 3\beta}{2 - 3\alpha\beta + 2i\omega\Delta t\alpha\beta} \left(u_n - 2\alpha\beta u_{n-1} + \frac{\alpha\beta}{2}u_{n-2} \right).$$

Then identifying the expression for v_{n+1} with the one obtained from v_n after shifting indices $n \rightarrow n+1$, we infer that the hoRAW filtered leapfrog scheme is equivalent to the following linear multistep method

$$u_{n+1} = \left(\frac{\alpha\beta + 3\beta}{2} \right) u_n + (1 - 2\beta)u_{n-1} - \left(\frac{\alpha\beta - \beta}{2} \right) u_{n-2}$$

$$+ i\omega\Delta t \left((2 - \beta + \alpha\beta)u_n - 3\alpha\beta u_{n-1} + \alpha\beta u_{n-2} \right). \quad (2.2)$$

Therefore the numerical amplification factor $A = \frac{u_{n+1}}{u_n}$ of the LF-hoRAW method satisfies the characteristic equation

$$A^3 - \left(\frac{\alpha\beta + 3\beta}{2} + (2 + \alpha\beta - \beta)i\omega\Delta t \right) A^2$$

$$- (1 - 2\beta - 3i\omega\Delta t\alpha\beta)A + \frac{\alpha\beta - \beta}{2} - i\omega\Delta t\alpha\beta = 0, \quad (2.3)$$

with one of the three roots, the physical mode, denoted by A_+ , and two computational modes.⁴ The exact solution $u(t) = \exp(i\omega t)u(0)$ of oscillation equation (2.1) has the exact amplification factor $A_{exact} = \exp(i\omega\Delta t)$. The behaviour of the exact and numerical amplification factors of LF-hoRAW scheme in the complex plane is shown in Figure 2.2 for various α and β . The exact amplification factor remains on the unit circle when $\omega\Delta t$ increases from 0 to 1. One of the computational modes of LF-hoRAW is amplified when $\omega\Delta t \geq \Sigma^{\alpha\beta}$ (see equation (2.4)) while the physical mode A_+ of LF-hoRAW stays inside the unit circle, similarly to the physical modes of AB3 [15, 16] and LF-hoRA [33]. The magnitudes of the physical and computational modes of LF-hoRAW are shown in Figure 2.3 for various values of α and β . This indicates that the LF-hoRAW scheme successfully controls the growth of its computational modes within the stability interval.

⁴The roots of (2.3) are obtained using Matlab's Symbolic Math Toolbox.

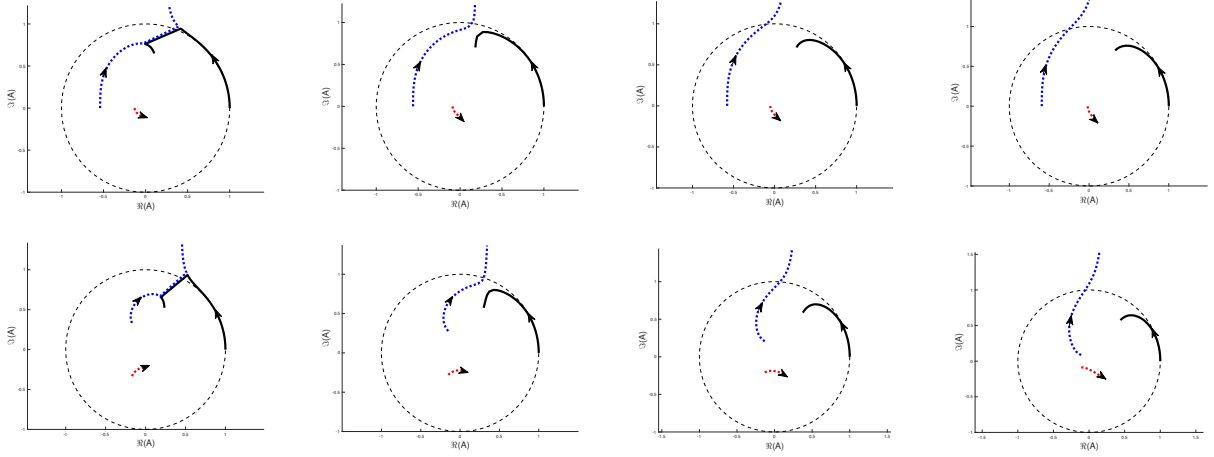


Figure 2.2: The amplification factors of the physical mode (solid line) and two computational modes (dotted line) of LF-hoRAW. From left to right: $\alpha = 0.3, 0.5, 0.7, 0.9$, with $\beta = 0.2$ (top) and $\beta = 0.4$ (bottom).

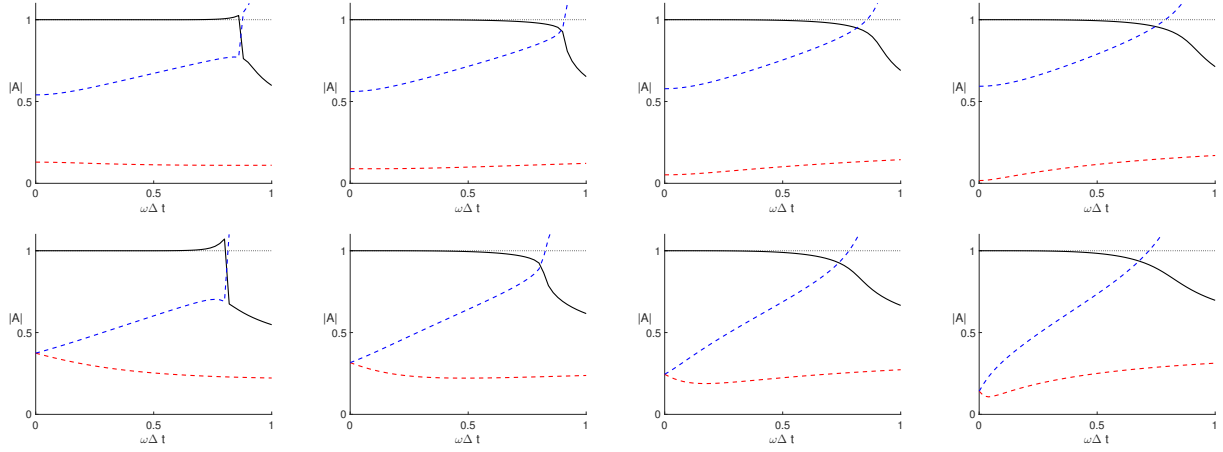


Figure 2.3: The magnitudes of the physical mode (solid line) and computational modes (dotted line) of LF-hoRAW. From left to right: $\alpha = 0.3, 0.5, 0.7, 0.9$, with $\beta = 0.2$ (top) and $\beta = 0.4$ (bottom).

2.1.3 The Consistency Order of LF-hoRAW

Using the Taylor expansions of $u(t_{n+1})$, $u(t_{n-1})$ and $u(t_{n-2})$ at time t_n , the local truncation error of LF-hoRAW (2.2) writes

$$\begin{aligned}\tau_{n+1}(\Delta t) &= \frac{1}{\Delta t} \left(u(t_{n+1}) - \frac{\alpha\beta + 3\beta}{2} u(t_n) - (1 - 2\beta) u(t_{n-1}) + \frac{\alpha\beta - \beta}{2} u(t_{n-2}) \right) \\ &\quad - (2 + \alpha\beta - \beta) i\omega u(t_n) + 3\alpha\beta i\omega u(t_{n-1}) - \alpha\beta i\omega u(t_{n-2}) \\ &= \frac{2 + 2\beta - 7\alpha\beta}{6} (i\omega\Delta t)^2 u'(t_n) + \frac{28\alpha\beta - 6\beta}{24} (i\omega\Delta t)^3 u'(t_n) + \mathcal{O}((i\omega\Delta t)^4).\end{aligned}$$

Therefore, the LF-hoRAW scheme exhibits third-order accuracy when $\alpha = \frac{2+2\beta}{7\beta}$, otherwise second-order.

2.1.4 The Stability Domain of LF-hoRAW

We determine the maximum interval of $\omega\Delta t$ for which all numerical amplification factors of LF-hoRAW scheme are non-amplified using the root locus curve method [25]. The characteristic equation of LF-hoRAW (2.2) is

$$\zeta^3 - \left(\frac{\alpha\beta + 3\beta}{2}\right)\zeta^2 - (1 - 2\beta)\zeta + \left(\frac{\alpha\beta - \beta}{2}\right) - z((2 + \alpha\beta - \beta)\zeta^2 - 3\alpha\beta\zeta + \alpha\beta) = 0$$

where $\zeta = \exp(i\theta)$, $\theta \in [0, 2\pi]$ represent the points on the unit circle, and $z \in \mathbb{C}$ is the root locus curve (see Figure 2.4). The stability interval of LF-hoRAW is determined by the intersection of the imaginary axis with the root locus curve z . Setting the real part of ζ to zero gives

$$\cos(\theta) = 1 \quad \text{or} \quad \cos(\theta) = \frac{5\alpha\beta - 4\alpha - \beta + 2}{4\alpha},$$

and

$$z = 0 \quad \text{or} \quad z = \pm i \frac{(2 + \alpha\beta - \beta)\sqrt{\beta + 8\alpha - 5\alpha\beta - 2}}{2\alpha(2 - \beta)\sqrt{2 + 5\alpha\beta - \beta}},$$

which are the intersections of the root locus curve with the imaginary axis. Therefore, LF-hoRAW is stable provided

$$\omega\Delta t \leq \Sigma^{\alpha\beta},$$

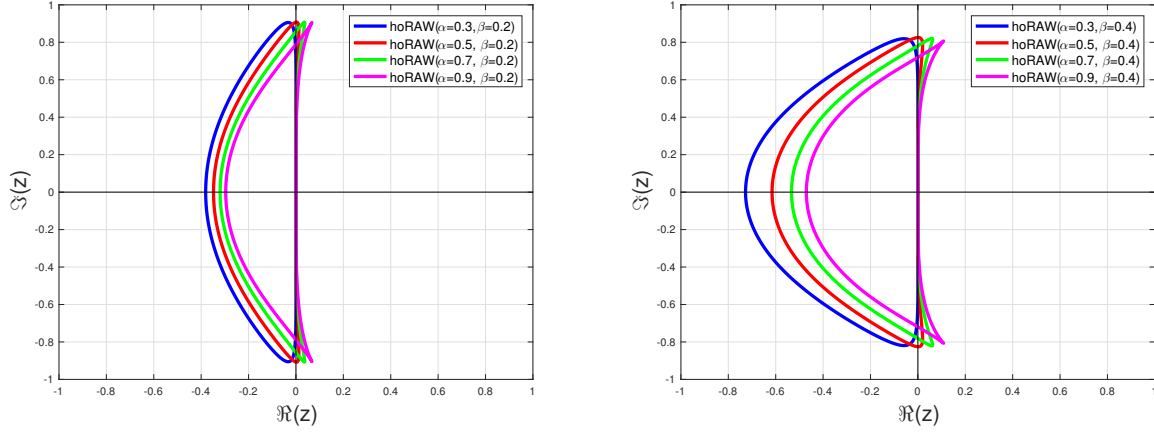


Figure 2.4: Root locus curve of LF-hoRAW with various α and β .

where

$$\Sigma^{\alpha\beta} = \frac{(2 + \alpha\beta - \beta)\sqrt{\beta + 8\alpha - 5\alpha\beta - 2}}{2\alpha(2 - \beta)\sqrt{2 + 5\alpha\beta - \beta}}, \quad (2.4)$$

with $\beta \in (0, 1)$ and $\alpha \in (0, 1]$. For any given $\beta \in (0, 1)$, the optimal value of α which maximizes $\Sigma^{\alpha\beta}$ is

$$\alpha^s = \frac{4 - 12\beta + 5\beta^2 - 2\sqrt{4 + 12\beta - 15\beta^2 + 4\beta^3}}{25\beta^2 - 36\beta}. \quad (2.5)$$

Hence, LF-hoRAW attains maximum stability when $\alpha = \alpha^s$, i.e.,

$$\omega\Delta t \leq \Sigma^{\alpha^s\beta} = \frac{\sqrt{2}(\sqrt{1 + 4\beta} + 1)^{\frac{1}{2}}(17 - 10\beta + \sqrt{1 + 4\beta})^{\frac{3}{2}}}{(2 - \beta)(13 + 5\sqrt{1 + 4\beta})^{\frac{3}{2}}}, \quad (2.6)$$

(see also Figure 2.9). From equation (2.4) we also note that the method becomes unstable when

$$\alpha = 1 - \frac{2}{\beta} \text{ or } \alpha = \frac{2 - \beta}{8 - 5\beta}.$$

Since for $\beta \in (0, 1)$ we have $1 - \frac{2}{\beta} < 0 < \frac{2 - \beta}{8 - 5\beta} < 1$, henceforth we will only consider $\alpha \in (\alpha^a, 1]$, where

$$\alpha^a = \frac{2 - \beta}{8 - 5\beta}. \quad (2.7)$$

We shall see in Section 2.3 that even if the scheme is unconditionally unstable when $\alpha = \alpha^a$, for slightly larger values $\alpha \gtrsim \alpha^a$, the LF-hoRAW is conditionally stable and the solution achieves almost sixth-order accuracy in amplitude. This phenomenon is similar to LF-RAW [48]. The amplitudes of the physical mode of LF-hoRAW are plotted in Figure 2.5 for several $\alpha \gtrsim \alpha^a$, and given $\beta = 0.2$ and 0.4 .

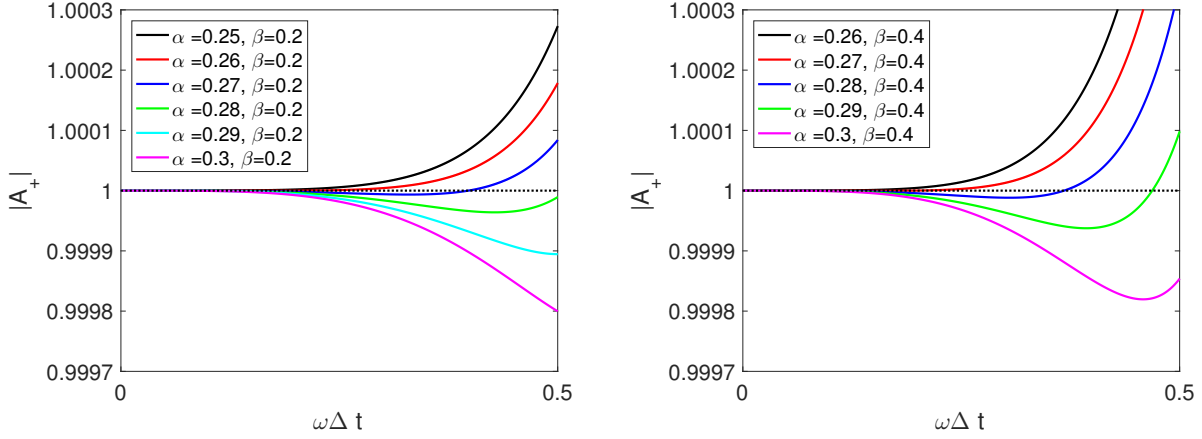


Figure 2.5: The magnitude of physical mode amplitudes while $\alpha \gtrsim \frac{2-\beta}{8-5\beta}$ for $\beta = 0.2$ (left) and $\beta = 0.4$ (right).

2.2 CURVATURE EVOLUTION

This section gives a geometric interpretation of the hoRAW filter in terms of the curvature evolution [28, 31]. We define the discrete curvature of φ^n by

$$\kappa(\varphi^n) = \varphi^{n+1} - 2\varphi^n + \varphi^{n-1}.$$

Two discrete curvatures are computed at every time integration of the system (1.4), one before and one after the time filter:

$$\kappa_{\text{old}}^n = w_{n+1} - 2v_n + u_{n-1}, \quad \kappa_{\text{new}}^n = v_{n+1} - 2u_n + u_{n-1}.$$

Figure 2.6 illustrates how the hoRAW time filter reduces the discrete curvature of the solution. After solving for w_{n+1} in the LF step (LF) the first solution curve is the continuous line. The curvature obtained is κ_{old}^n . Next, performing the (hoRA) and (W) steps leads to the new solution curve (the dashed line of Figure 2.6), with curvature κ_{new}^n .

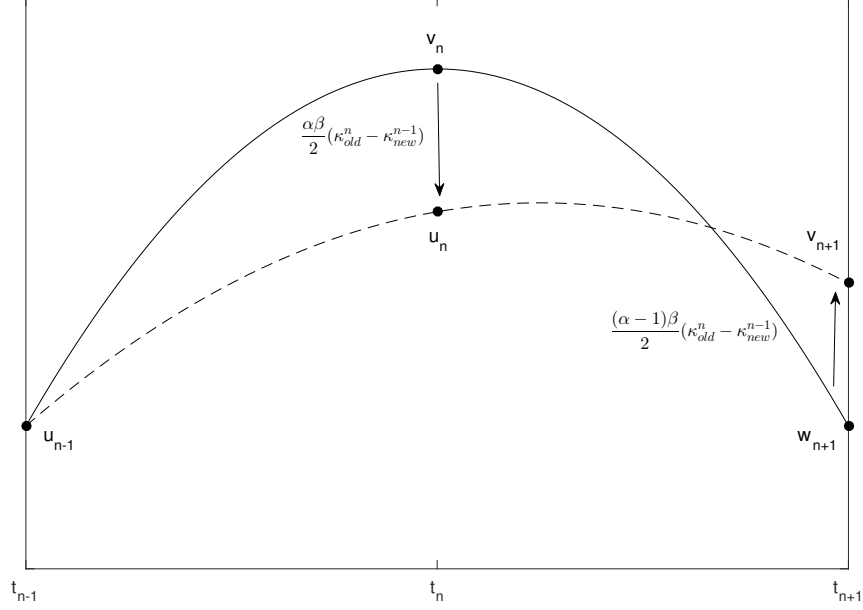


Figure 2.6: The hoRAW filter moves the inner and right outer points through displacements $\alpha\beta(\kappa_{\text{old}}^n - \kappa_{\text{new}}^{n-1})/2$ and $(1-\alpha)\beta(\kappa_{\text{old}}^n - \kappa_{\text{new}}^{n-1})/2$, respectively, where $\alpha \in (\alpha^a, 1]$, $\beta \in [0, 1]$. The standard hoRA filter moves only the inner point through a displacement $\beta(\kappa_{\text{old}}^n - \kappa_{\text{new}}^{n-1})/2$.

The next result shows that the hoRAW filter preserves the three-time-level mean of the solution curve, and decreases the discrete curvature of the solution.

Proposition 2.2.1. *For $n \geq 1$, we have*

$$\kappa_{\text{new}}^n \leq \max\{\kappa_{\text{new}}^{n-1}, \kappa_{\text{old}}^n\}$$

for all $\alpha \in [\alpha^a, 1]$, $\beta \in [0, 1]$. When $\alpha = 1/2$, the hoRAW filter preserves the three-time-level mean of the solution curves

$$\frac{v_{n+1} + u_n + u_{n-1}}{3} = \frac{w_{n+1} + v_n + u_{n-1}}{3}.$$

Proof. The first step (hoRA) of the hoRAW filter is $u_n = v_n + \frac{\beta\alpha}{2}\kappa_{\text{old}}^n - \frac{\beta\alpha}{2}\kappa_{\text{new}}^{n-1}$ and the second step (W) of the hoRAW filter is $v_{n+1} = w_{n+1} + \frac{\beta(\alpha-1)}{2}\kappa_{\text{old}}^n - \frac{\beta(\alpha-1)}{2}\kappa_{\text{new}}^{n-1}$. Then

$$\begin{aligned}\kappa_{\text{new}}^n &= v_{n+1} - 2u_n + u_{n-1} \\ &= (v_{n+1} - w_{n+1}) + 2(v_n - u_n) + (w_{n+1} - 2v_n + u_{n-1}) \\ &= \frac{\beta(\alpha-1)}{2}(\kappa_{\text{old}}^n - \kappa_{\text{new}}^{n-1}) + \beta\alpha(\kappa_{\text{new}}^{n-1} - \kappa_{\text{old}}^n) + \kappa_{\text{old}}^n \\ &= \frac{\beta(\alpha+1)}{2}\kappa_{\text{new}}^{n-1} + \left(1 - \frac{\beta(\alpha+1)}{2}\right)\kappa_{\text{old}}^n.\end{aligned}$$

The claim follows by taking the maximum of $\{\kappa_{\text{new}}^{n-1}, \kappa_{\text{old}}^n\}$. Adding the (hoRA) and (W) steps of the hoRAW filter, for $\alpha = 1/2$, yields

$$\begin{aligned}\frac{v_{n+1} + u_n + u_{n-1}}{3} &= \frac{w_{n+1} + v_n + \frac{\beta(2\alpha-1)}{2}\kappa_{\text{old}}^n - \frac{\beta(2\alpha-1)}{2}\kappa_{\text{new}}^{n-1} + u_{n-1}}{3} \\ &= \frac{w_{n+1} + v_n + u_{n-1}}{3},\end{aligned}$$

hence the three-time-level means are preserved. \square

2.3 ERROR ANALYSIS FOR PHASE AND AMPLITUDE

We will derive the phase and amplitude errors of the LF-hoRAW scheme (2.2) using modified equation [5, 20, 26, 46].

We define the following real constants $C_1^{\alpha\beta}$, $C_2^{\alpha\beta}$ and $C_3^{\alpha\beta}$, depending on $\alpha \in [\alpha^a, 1]$, $\beta \in [0, 1]$ as:

$$\begin{aligned}C_1^{\alpha\beta} &= \frac{2 + 2\beta - 7\alpha\beta}{6(2 - \beta - \alpha\beta)}, \\ C_2^{\alpha\beta} &= \frac{5\alpha\beta^2 - 8\alpha\beta + 2\beta - \beta^2}{4(2 - \beta - \alpha\beta)^2}, \\ C_3^{\alpha\beta} &= \frac{\beta^3(113\alpha^3 + 54\alpha^2 - 101\alpha + 18) - 2\beta^2(169\alpha^2 - 112\alpha + 9) + 4\beta(19\alpha - 6) - 24}{40(\alpha\beta + \beta - 2)^3},\end{aligned}\tag{2.8}$$

and the three term modified equation of LF-hoRAW corresponding to (2.1):

$$x'(t) = \left(1 - C_1^{\alpha\beta}(i\omega\Delta t)^2 + C_2^{\alpha\beta}(i\omega\Delta t)^3 + C_3^{\alpha\beta}(i\omega\Delta t)^4\right) i\omega x(t).\tag{2.9}$$

Proposition 2.3.1. *The LF-hoRAW (2.2) is a fifth-order approximation to the modified equation (2.9)*

$$\widehat{\tau}(\Delta t) = G(x)\Delta t^5,$$

while only second-order approximation to the oscillation equation (2.1).

Proof. Consider a general three term modified equation corresponding to the oscillation equation

$$y'(t) = i\omega y(t) + \Delta t^2 g_1(y(t)) + \Delta t^3 g_2(y(t)) + \Delta t^4 g_3(y(t)). \quad (2.10)$$

Then the local truncation error of LF-hoRAW (based on the modified equation, not on the oscillation equation) is

$$\begin{aligned} \widehat{\tau}_{n+1}(\Delta t) = & \frac{1}{\Delta t} \left[y(t_{n+1}) - \left(\frac{\alpha\beta + 3\beta}{2} \right) y(t_n) - (1 - 2\beta) y(t_{n-1}) + \left(\frac{\alpha\beta - \beta}{2} \right) y(t_{n-2}) \right] \\ & - (2 + \alpha\beta - \beta) i\omega y(t_n) + 3\alpha\beta i\omega y(t_{n-1}) - \alpha\beta i\omega y(t_{n-2}). \end{aligned}$$

Using the Taylor expansions of $y(t_{n+1})$, $y(t_{n-1})$, $y(t_{n-2})$ at time t_n , and substitute in $y^{(i)}(t_n)$, $i = 1, \dots, 5$, the local truncation error writes

$$\begin{aligned} \widehat{\tau}_{n+1}(\Delta t) = & \left((2 - \beta - \alpha\beta) g_1(y(t_n)) + \frac{\alpha\beta}{2} i\omega^3 y(t_n) - \frac{2 - 4\alpha\beta + 2\beta}{6} i\omega^3 y(t_n) \right) \Delta t^2 \\ & + \left((2 - \beta - \alpha\beta) g_2(y(t_n)) + \alpha\beta i\omega g_1'(y(t_n)) y(t_n) + \frac{28\alpha\beta - 6\beta}{24} \omega^4 y(t_n) \right) \Delta t^3 \\ & + \left((2 - \beta - \alpha\beta) g_3(y(t_n)) + \alpha\beta i\omega g_2'(y(t_n)) y(t_n) + \frac{\alpha\beta}{2} \omega^2 g_1(y(t_n)) \right. \\ & + \frac{\alpha\beta}{2} \omega^2 g_1'(y(t_n)) y(t_n) - \frac{2 - 4\alpha\beta + 2\beta}{6} \omega^2 g_1(y(t_n)) \\ & \left. - \frac{4 - 8\alpha\beta + 4\beta}{6} \omega^2 g_1'(y(t_n)) y(t_n) + \frac{2 - 81\alpha\beta + 14\beta}{120} i\omega^5 y(t_n) \right) \Delta t^4 + \mathcal{O}(\Delta t^5). \end{aligned}$$

Setting the coefficients of Δt^2 , Δt^3 and Δt^4 to zero, we obtain that $g_1(y) = C_1^{\alpha\beta} i\omega^3 y$, $g_2(y) = C_2^{\alpha\beta} \omega^4 y$ and $g_3(y) = C_3^{\alpha\beta} i\omega^5 y$, concluding the proof. \square

Recall that the global error based on the modified equation (2.9) is $x(t_n) - u_n$ and it coincides with the truncation error $\Delta t \hat{\tau}_n(\Delta t)$ under the localizing assumption $u_{n-i} = x(t_{n-i})$, $i = 1, 2, 3$. Thus, we have

$$x(t_n) - u_n = \mathcal{O}(\Delta t^6),$$

by Proposition 2.3.1. Therefore, the global error of LF-hoRAW

$$u(t_n) - u_n = u(t_n) - x(t_n) + x(t_n) - u_n = u(t_n) - x(t_n) + \mathcal{O}(\Delta t^6)$$

can be characterized by the difference between the curves $u(t)$ and $x(t)$.

Theorem 2.3.1. *The phase and amplitude errors of the LF-hoRAW scheme applied to oscillation equation (2.1) are*

$$R_+ - 1 = \frac{2 + 2\beta - 7\alpha\beta}{6(2 - \beta - \alpha\beta)}(\omega\Delta t)^2 + \left[\frac{\beta^3(113\alpha^3 + 54\alpha^2 - 101\alpha + 18)}{40(\alpha\beta + \beta - 2)^3} - \frac{2\beta^2(169\alpha^2 - 112\alpha + 9) - 4\beta(19\alpha - 6) + 24}{40(\alpha\beta + \beta - 2)^3} \right](\omega\Delta t)^4 + \mathcal{O}((\omega\Delta t)^6), \quad (2.11)$$

$$|A_+| - 1 = \frac{5\alpha\beta^2 - 8\alpha\beta + 2\beta - \beta^2}{4(2 - \beta - \alpha\beta)^2}(\omega\Delta t)^4 + \mathcal{O}((\omega\Delta t)^6).$$

Proof. With the initial conditions $u(0) = x(0) = 1$, the exact solution of oscillation equation (2.1) is $u(t) = \exp(i\omega t)$ and the exact solution to the modified equation (2.9) is

$$\begin{aligned} x(t) &= \exp(i\omega t + C_1^{\alpha\beta}i\omega^3(\Delta t)^2t + C_2^{\alpha\beta}\omega^4(\Delta t)^3t + C_3^{\alpha\beta}i\omega^5(\Delta t)^4t) \\ &= \exp(C_2^{\alpha\beta}\omega^4(\Delta t)^3t) \left(\cos(\omega t + C_1^{\alpha\beta}\omega^3(\Delta t)^2t + C_3^{\alpha\beta}\omega^5(\Delta t)^4t) \right. \\ &\quad \left. + i \sin(\omega t + C_1^{\alpha\beta}\omega^3(\Delta t)^2t + C_3^{\alpha\beta}\omega^5(\Delta t)^4t) \right) \end{aligned}$$

where $C_1^{\alpha\beta}$, $C_2^{\alpha\beta}$ and $C_3^{\alpha\beta}$ are defined in (2.8). Thus, the phase and amplitude errors of the LF-hoRAW method in one time step are

$$R_+ - 1 = \frac{\arg(x(\Delta t))}{\arg(u(\Delta t))} - 1 = C_1^{\alpha\beta}(\omega\Delta t)^2 + C_3^{\alpha\beta}(\omega\Delta t)^4 + \mathcal{O}((\omega\Delta t)^6),$$

$$|A_+| - 1 = |x(\Delta t)| - |u(\Delta t)| = |\exp(C_2^{\alpha\beta}(\omega\Delta t)^4)| - 1 = C_2^{\alpha\beta}(\omega\Delta t)^4 + \mathcal{O}((\omega\Delta t)^6),$$

in the limit of good time resolution $\omega\Delta t \ll 1$. □

Remark 2.3.1. The phase and amplitude errors of the LF-hoRA [33] are recovered when $\alpha = 1$, and also the phase and amplitude errors of the leapfrog method are recovered when $\alpha = 1, \beta = 0$.

Remark 2.3.2. The LF-hoRAW method attains sixth-order accuracy in amplitude when $\alpha = \alpha^a$ and fourth-order accuracy in phase when $\alpha = \frac{2+2\beta}{7\beta}$.

The amplitude and the relative phase change of the physical mode of LF-hoRAW are illustrated in Figure 2.7 with various of α and β . Next, the summary of stability, accuracy in terms of phase and amplitude, and conservation of three-time-level mean for LF-hoRAW is presented in Table 2.1.

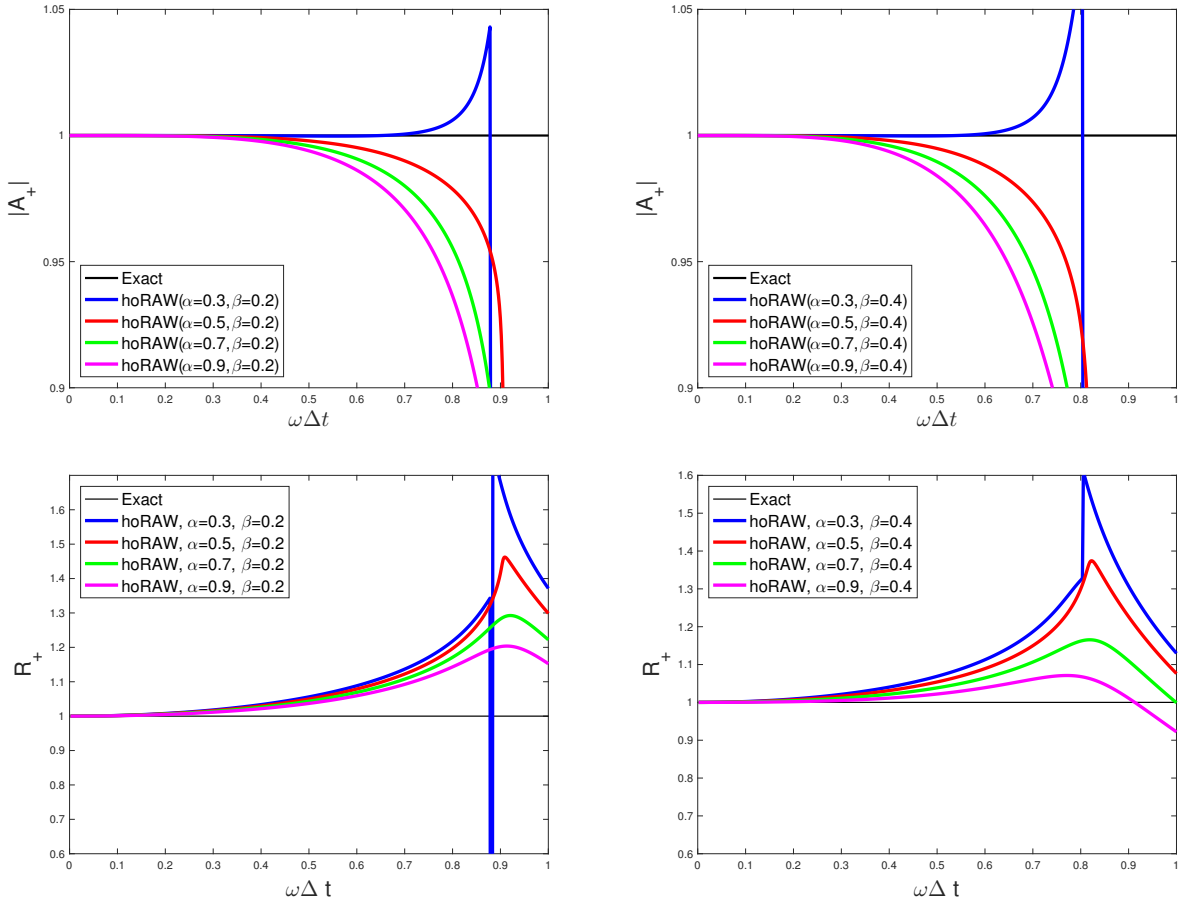


Figure 2.7: Amplitude (top) and relative phase change (bottom) of the physical mode of LF-hoRAW for $\alpha = 0.3, 0.5, 0.7, 0.9$ with $\beta = 0.2$ (left) and $\beta = 0.4$ (right).

Table 2.1: Summary of the conservation of three-time-level mean, stability, and accuracy properties of the LF-hoRAW for some values of α .

α	Conserves three		Order of accuracy	
	time level mean	Stability	Amplitude	Phase [*]
$= \alpha^a(2.7)$	No	Unconditionally Unstable	6	2 (resp. 4)
$\gtrsim \alpha^a(2.7)$	No	Conditionally Stable	4	2 (resp. 4)
1/2	Yes	Conditionally Stable	4	2 (resp. 4)
1	No	Conditionally Stable	4	2 (resp. 4)

^{*} The phase is fourth-order when $\alpha = \frac{2+2\beta}{7\beta}$.

2.4 COMPARISON OF LF-HORA, LF-HORAW AND AB3 METHODS

First we summarize the properties of the third-order methods LF-hoRA ($\beta = 0.4$) [33], LF-hoRAW ($\alpha = \frac{2+2\beta}{7\beta}$, $\beta \in [0, 1]$) and AB3 [15].

- The truncation errors are:

$$\tau_n(\Delta t) = \frac{11}{30}(i\omega\Delta t)^3 u'(t_n) + \mathcal{O}((i\omega\Delta t)^4) \quad (\text{hoRA}, \beta = 0.4)$$

$$\tau_n(\Delta t) = \frac{\beta+4}{12}(i\omega\Delta t)^3 u'(t_n) + \mathcal{O}((i\omega\Delta t)^4) \quad (\text{hoRAW}, \alpha = \frac{2+2\beta}{7\beta})$$

$$\tau_n(\Delta t) = \frac{3}{8}(i\omega\Delta t)^3 u'(t_n) + \mathcal{O}((i\omega\Delta t)^4) \quad (\text{AB3})$$

When $\beta \in [0, 0.4]$, we have that $\frac{\beta+4}{12} \in [\frac{1}{3}, \frac{11}{30}]$, therefore LF-hoRAW method has a smaller truncation error compared to LF-hoRA and AB3.

- The methods are stable for:

$$\omega\Delta t \leq 0.69 \quad (\text{hoRA}, \beta = 0.4)$$

$$\omega\Delta t \leq \frac{(16-5\beta)\sqrt{\beta}\sqrt{16-8\beta-3\beta^2}}{4\sqrt{3}(2-\beta)(1+\beta)\sqrt{\beta+8}} \quad (\text{hoRAW}, \alpha = \frac{2+2\beta}{7\beta})$$

$$\omega\Delta t \leq 0.72 \quad (\text{AB3})$$

In long-time integrations, probably the most important quantities are the stability interval and the amplitude accuracy. The comparison of stability regions determined by root locus curves for LF-hoRAW, LF-hoRA and AB3 methods are plotted in Figure 2.8. The comparison of stability and leading coefficient of amplitude error are shown in Figure 2.9.

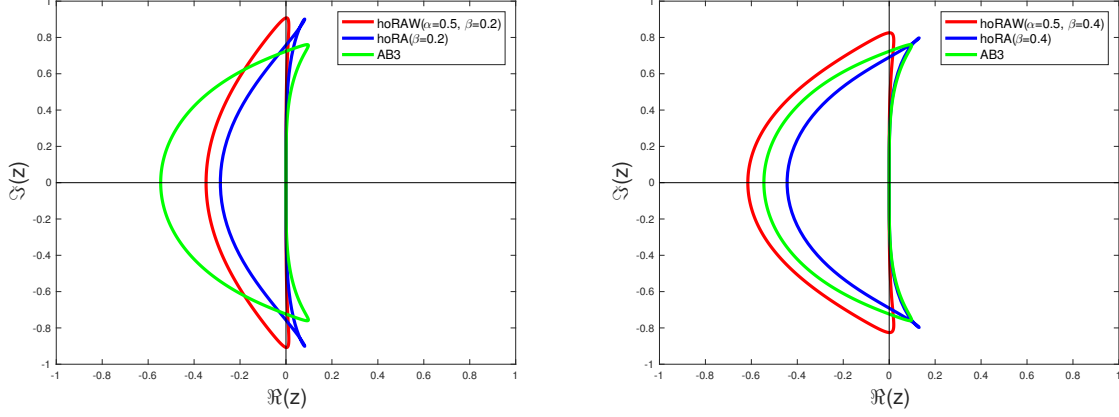


Figure 2.8: Root locus curves of LF-hoRA, LF-hoRAW and AB3 schemes. The stability interval is given by the intersection of the root locus curve with the imaginary axis.

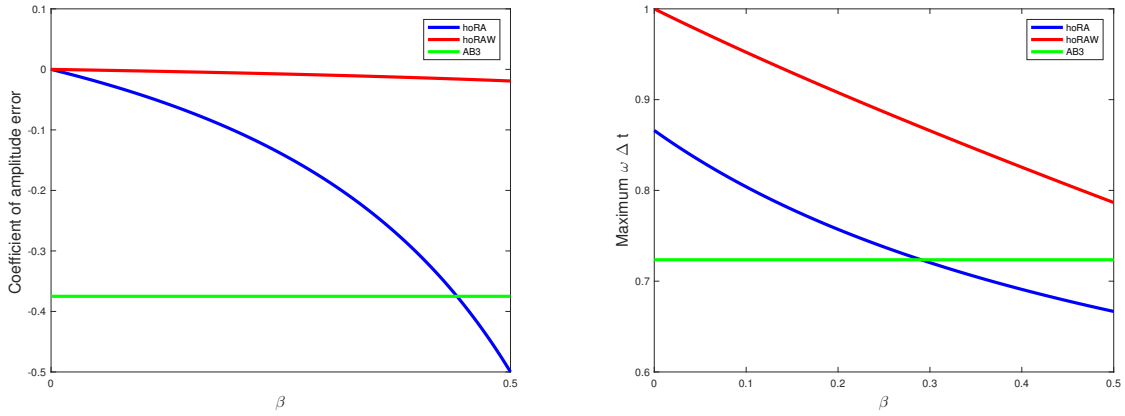


Figure 2.9: Comparison of the coefficients $C_2^{1\beta}$, $C_2^{\alpha\beta}$ (2.8) and $C_{AB3} = 3/8$ corresponding the LF-hoRA, LF-hoRAW and AB3 amplitude errors (2.11) (left). Comparison of the maximum stability for LF-hoRA ($\Sigma^{1\beta}$), LF-hoRAW (*see*(2.6)) and AB3 (0.7236) (right).

The phase and amplitude of the physical mode of LF-hoRAW, LF-hoRA and AB3 are illustrated in Figure 2.10. We next compare the accuracy of amplitude and the stability intervals of the second-order LF-hoRAW with the second and third-order LF-hoRA, and the AB3 methods in Table 2.2, with some featured values of α and β .

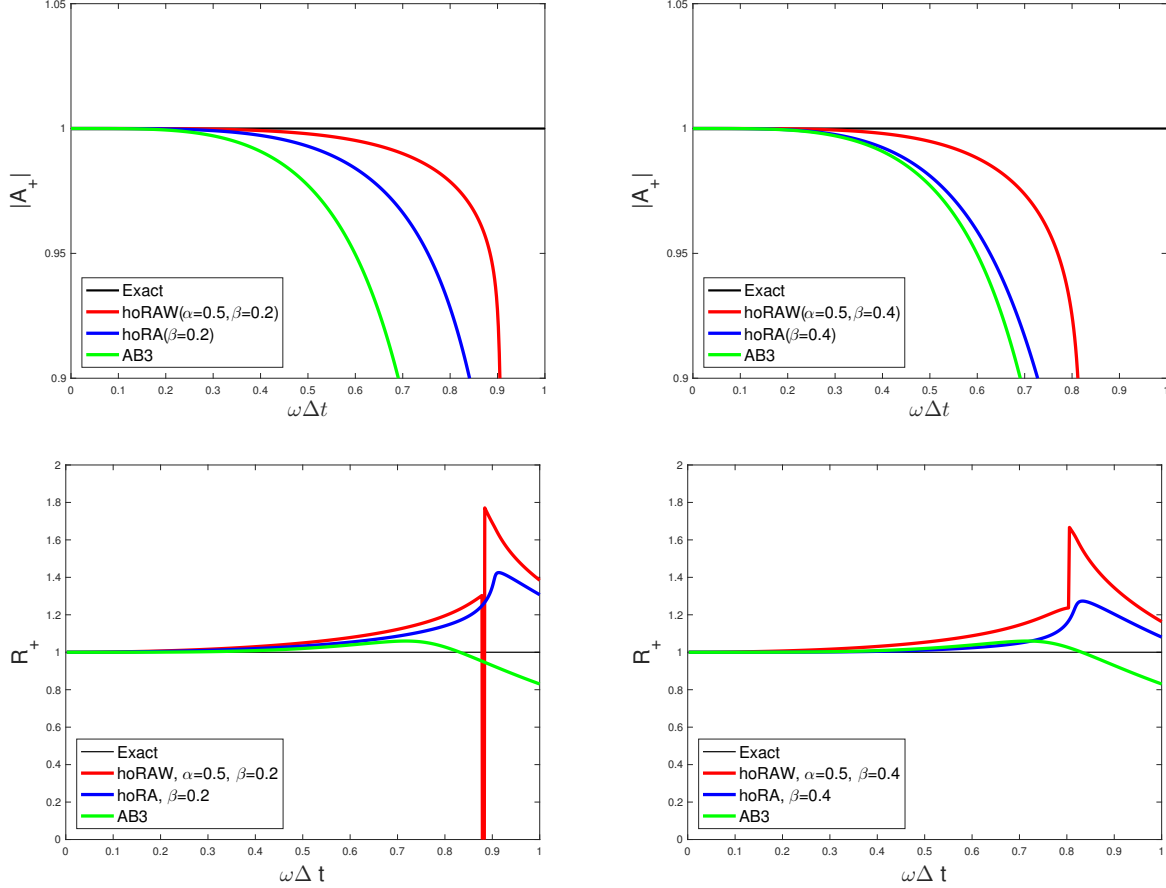


Figure 2.10: The comparison of amplitude (top) and relative phase (bottom) of physical mode of LF-hoRA, LF-hoRAW and AB3.

From the Table 2.2 and Figure 2.9 we infer that the LF-hoRAW ($\beta = 0.2, \alpha = \alpha^s(\beta) \equiv 0.4887$) method exhibits about a twenty percent increase in stability compared to LF-hoRA method, and twenty five percent compared to the three-level intrusive AB3 method. Moreover, the LF-hoRAW increases the amplitude accuracy by one significant digit. Also, when $\beta = 0.2, \alpha = 0.27 \gtrsim \alpha^a(0.2) = 0.2571$ and $\beta = 0.4, \alpha = 0.28 \gtrsim \alpha^a(0.2) = 0.2667$ the increase in amplitude accuracy is by two significant digits.

Table 2.2: The comparison of LF-hoRAW, LF-hoRA and AB3 schemes with some featured values of α and β

Method	β	α	Order	Amplitude	Max. $\omega\Delta t$
LF	-	-	2	1	1
LF-hoRA	0.2	-	2	$1 - .1016(\omega\Delta t)^4$	0.7571
LF-hoRAW	0.2	0.27	2	$1 - .0015(\omega\Delta t)^4$	0.3977
LF-hoRAW	0.2	0.3	2	$1 - .0050(\omega\Delta t)^4$	0.6509
LF-hoRAW	0.2	0.4887	2	$1 - .0280(\omega\Delta t)^4$	0.9078
LF-hoRAW	0.2	0.5	2	$1 - .0294(\omega\Delta t)^4$	0.9075
LF-hoRA	0.4	-	3	$1 - .3056(\omega\Delta t)^4$	0.6910
LF-hoRAW	0.4	0.28	2	$1 - .0036(\omega\Delta t)^4$	0.3677
LF-hoRAW	0.4	0.3	2	$1 - .0091(\omega\Delta t)^4$	0.5402
LF-hoRAW	0.4	0.4961	2	$1 - .0701(\omega\Delta t)^4$	0.8256
LF-hoRAW	0.4	0.5	2	$1 - .0714(\omega\Delta t)^4$	0.8255
AB3	-	-	3	$1 - .3750(\omega\Delta t)^4$	0.7236

2.5 NUMERICAL TESTS

We present three numerical tests on the LF-hoRAW, LF-hoRA and the AB3 methods to validate the stability and error analysis presented in previous section.

2.5.1 Simple Pendulum

Consider a simple pendulum problem, given by the following two coupled nonlinear equations (see [33, 50])

$$\begin{aligned}
 \frac{d\theta}{dt} &= \frac{v}{L}, \\
 \frac{dv}{dt} &= -g \sin \theta,
 \end{aligned}
 \tag{2.12}$$

where θ, v, L and g denote, respectively, the angular displacement, velocity along the arc, length of the pendulum, and the acceleration due to gravity.

Set the initial condition $(\theta(0), v(0)) = (0.9\pi, 0)$ close to the unstable equilibrium point, $g = 9.8$, $L = 49$, and numerically integrate the third-order LF-hoRA ($\alpha = 1, \beta = 0.4$), the third-order AB3, and the second-order LF-hoRAW ($\alpha = 0.3, \beta = 0.4$) over the time interval $[0, 400]$, with the time step $\Delta t = 0.5$. The fourth-order Runge-Kutta method is used to initialize the second and third steps for hoRA, hoRAW and AB3. We then compare the results with the reference solution, which computed using the adaptive RK4-5 method with the relative error tolerance 10^{-10} and absolute error tolerance 10^{-15} . The comparison is shown in Figure 2.11. The plots show that the LF-hoRA and AB3 damp the amplitude, and the phase errors are relatively large. The hoRAW filter preserves the phase and amplitude, with high accuracy, for a long-time period.

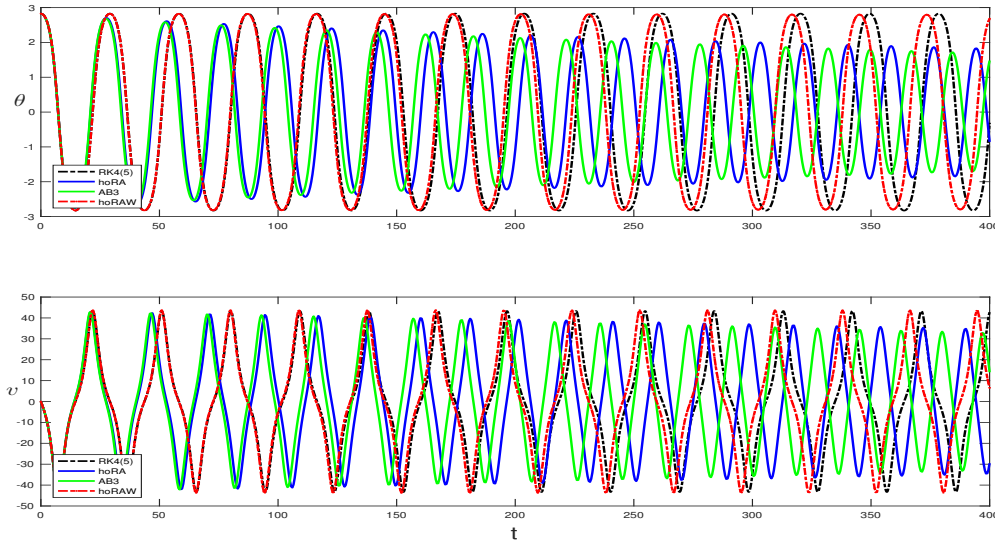
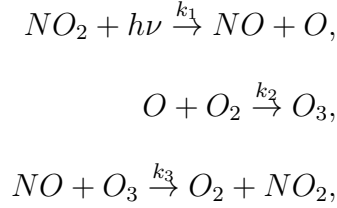


Figure 2.11: Numerical solution to the simple pendulum problem, computed by LF-hoRA ($\beta = 0.4$), AB3 and LF-hoRAW ($\alpha = 0.3, \beta = 0.4$) methods, are compared with the reference solution of the adaptive RK4-5 method with relative error tolerance 10^{-10} and absolute error tolerance 10^{-15} . The initial condition is $(\theta(0), v(0)) = (0.9\pi, 0)$, and the time step is $\Delta t = 0.5$.

2.5.2 Ozone Photochemistry

We consider a classic example from chemical reaction (see e.g., [27, 33]) between atomic oxygen(O), nitrogen oxides (NO and NO₂), and ozone (O₃):



where $h\nu$ denotes a photon of solar radiation. Assuming that the background concentration of O_2 is constant, the concentration $c = (c_1, c_2, c_3, c_4)$, in molecules per cubic centimeter, of O, NO, NO₂ and O₃, modeling the chemical reactions above, satisfies the system:

$$\begin{aligned} \frac{dc_1}{dt} &= k_1 c_3 - k_2 c_1, \\ \frac{dc_2}{dt} &= k_1 c_3 - k_3 c_2 c_4, \\ \frac{dc_3}{dt} &= k_3 c_2 c_4 - k_1 c_3, \\ \frac{dc_4}{dt} &= k_2 c_1 - k_3 c_2 c_4. \end{aligned}$$

We choose⁵

$$\begin{aligned} k_1 &= 10^{-2} \max\{0, \sin(2\pi t/t_d)\} s^{-1}, \\ k_2 &= 10^{-2} s^{-1}, \quad k_3 = 10^{-16} \text{cm}^3 \text{molecule}^{-1} s^{-1}, \end{aligned}$$

as in [33], where t_d is the length of 1 day in seconds, and the initial condition $c_0 = (0, 0, 5 \times 10^{11}, 8 \times 10^{11})$ molecules cm^{-3} at $t = 0$. The reference solution is computed using the adaptive RK4-5 method, with the same error tolerances as before. We compare two numerical solutions, LF-hoRA($\beta = 0.4$) and LF-hoRAW($\alpha = 0.3, \beta = 0.4$), computed with the time step $\Delta t = 45$ second. The chemical concentrations over the next 48 hours are shown in Figure 2.12. The LF-hoRAW method is able to capture the behaviour of concentrations with reasonable accuracy.

⁵ k_2 is chosen to make chemical reaction equations non-stiff instead of 10^5 as in [33]

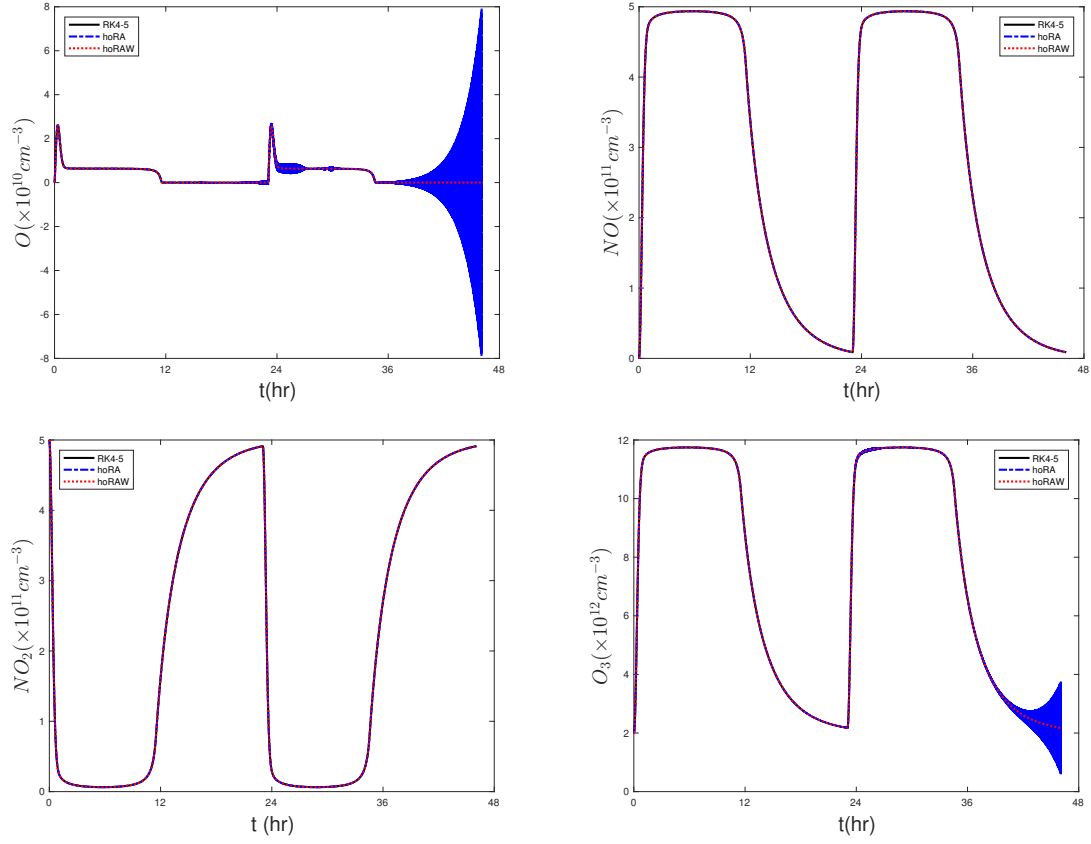


Figure 2.12: Numerical solutions for chemical concentrations, computed using LF-hoRAW ($\alpha = 0.3, \beta = 0.4$) and LF-hoRA($\beta = 0.4$), are compared with the reference solutions obtained from adaptive RK4-5 method with relative error tolerance 10^{-10} and absolute tolerance 10^{-15} . The initial condition is $c_0 = (0, 0, 5 \times 10^{11}, 8 \times 10^{11})$ molecules cm^{-3} at $t = 0$ with time step $\Delta t = 45$ second.

2.5.3 Lorenz System

Consider the Lorenz system:

$$\begin{aligned}
 \frac{dX}{dt} &= \sigma(Y - X), \\
 \frac{dY}{dt} &= -XZ + rX - Y, \\
 \frac{dZ}{dt} &= XY - bZ,
 \end{aligned} \tag{2.13}$$

with $\sigma = 12, r = 12, b = 6$, and the initial condition $(X_0, Y_0, Z_0) = (-10, -10, 25)$, as in [15, 33]. The system is numerically integrated over the time interval $[0, 5]$ using the third-order LF-hoRAW ($\alpha = 34/49, \beta = 0.7$), the third-order LF-hoRA ($\beta = 0.4$) and the third-order AB3 methods. The reference solution is computed with the adaptive RK4-5 method, using the same error tolerance as in 2.5.1. The numerical solutions of X are plotted in Figure 2.13 with various time steps $\Delta t = 0.025, 0.029, 0.035$ and 0.045 . The augmented stability property of LF-hoRAW is exhibited as the time step is increased, pushing the LF-hoRA and AB3 to become unstable and it also shows that the third-order LF-hoRAW is as accurate the third-order LF-hoRA and AB3 methods.

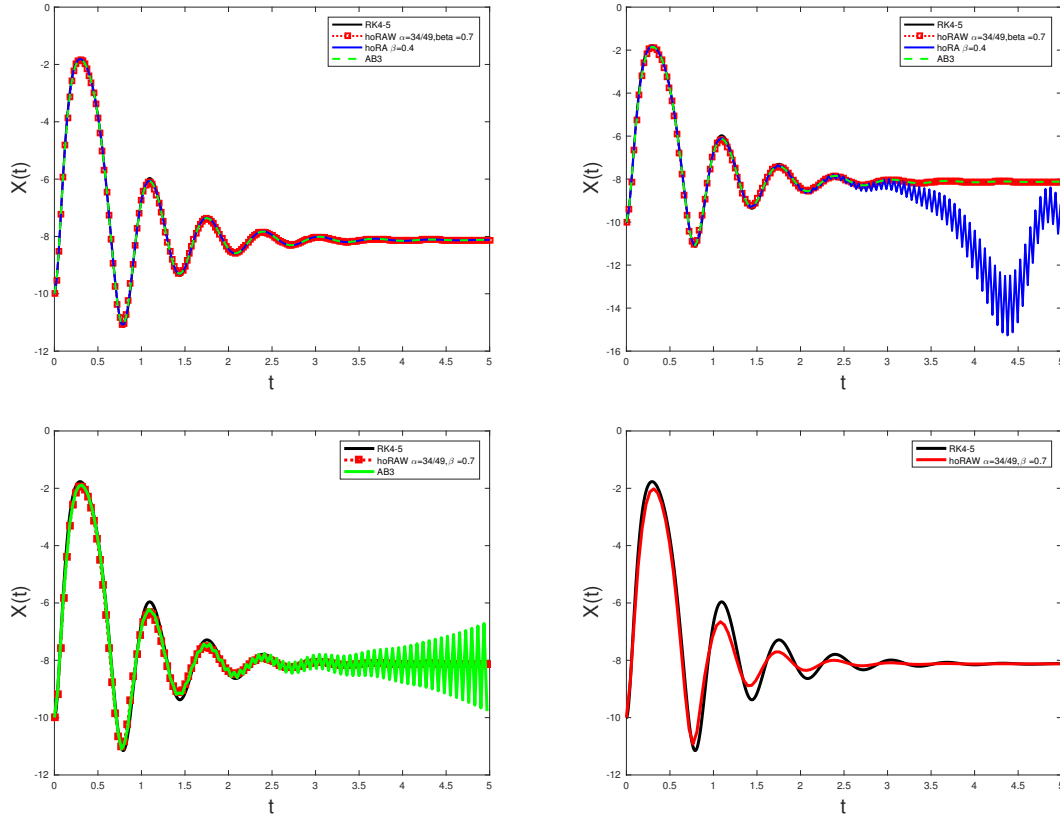


Figure 2.13: Computed numerical solutions to the Lorenz system with LF-hoRA($\beta = 0.4$), LF-hoRAW($\alpha = 34/49, \beta = 0.7$) and AB3 with constant time step $\Delta t = 0.025$ (top left), $\Delta t = 0.029$ (top right) and $\Delta t = 0.035$ (bottom left), $\Delta t = 0.045$ (bottom right).

2.6 SUMMARY

We have constructed and analyzed a higher order Robert-Asselin-Williams time filter. The LF solution is once-filtered in the (hoRA) step and twice-filtered in the (W) step. The LF-hoRAW has an increased stability and accuracy of amplitude when compared to LF-RA, LF-RAW, LF-hoRA and AB3 (see Figure 2.1, and Table 2.2, Figure 2.9). The effect of the twice-filtering with the Williams' step (W) is the increase by almost twenty percent in stability, compared to the only once-filtered LF-hoRA. The LF-hoRAW recovers the full-stability of LF when $\alpha = \alpha^s$, $\beta = 0$ unlike the hoRA filtered LF method (see Figure 2.9). The LF-hoRAW can achieve almost sixth-order accuracy in amplitude when $\alpha \gtrapprox \alpha^a(\beta)$ (see Theorem 2.3.1). The increase in amplitude accuracy is by two significant digits for the specific values $\beta = 0.2$, $\alpha = 0.27 \gtrapprox \alpha^a(0.2) = 0.2571$ and $\beta = 0.4$, $\alpha = 0.28 \gtrapprox \alpha^a(0.2) = 0.2667$. The hoRAW is an efficient and accurate post-process which successfully controls the growth of the computational modes and may be suitable for long-time simulations of weather and climate models.

3.0 THE GENERAL HIGH-ORDER ROBERT-ASSELIN TIME FILTER

In this chapter, we derive phase and amplitude error for general high-order Robert-Asselin (ghoRA) time filter proposed by Li in his PhD thesis [32]. We first briefly recall the properties of ghoRA.

3.1 PREVIOUS WORK

The ghoRA filtered leapfrog scheme applied to oscillation equation (2.1) is given by

$$v_{n+1} = u_{n-1} + 2i\omega\Delta t v_n \quad (\text{LF})$$

$$u_n = v_n + a_{n+1}v_{n+1} + a_n v_n + \sum_{j=1}^k a_{n-j}u_{n-j}, \quad k \geq 1. \quad (\text{ghoRA})$$

where a_j , $j = n+1, n, \dots, n-k$ are coefficient to be determined by Taylor expansion and v and u are unfiltered and once filtered values, respectively.

Theorem 3.1.1. *The general high-order Robert-Asselin (ghoRA) filtered leapfrog (LF) scheme (LF-ghoRA) is equivalent to following linear multistep method*

$$u_{n+1} - a_{n+1}u_n - (1 + a_n)u_{n-1} - \sum_{j=1}^k a_{n-j}u_{n+1-j} = 2i\omega\Delta t \left(u_n - \sum_{j=1}^k a_{n-j}u_{n-j} \right) \quad (3.1)$$

for all $k = 1, 2, \dots, n$.

Proof. Solving (LF) and (ghoRA) for v_n , we get

$$v_n = \frac{u_n - a_{n+1}u_{n-1} - \sum_{j=1}^k a_{n-j}u_{n-j}}{1 + 2i\omega\Delta t a_{n+1} + a_n}. \quad (3.2)$$

Similarly, solving (LF) and (ghoRA) for v_{n+1} , we obtain

$$v_{n+1} = \frac{(1 + a_n)u_{n-1} + 2i\omega\Delta t u_n - 2i\omega\Delta t \sum_{j=1}^k a_{n-j}u_{n-j}}{1 + a_n + 2i\omega\Delta t a_{n+1}} \quad (3.3)$$

Take $n \rightarrow n + 1$ in (3.2) and set to (3.3) gives the (3.1). \square

Theorem 3.1.2. *The linear multistep method (3.1) is consistent if and only if the following two condition holds,*

$$\sum_{j=-1}^k a_{n-j} = 0, \quad \sum_{j=0}^k (j+1)a_{n-j} = 0. \quad (3.4)$$

Furthermore, the linear multistep method (3.1) is p^{th} -order accurate for any $1 < p \leq k+1$ if and only if the condition (3.4) holds and following condition for $\ell = 2, \dots, p$ are satisfied,

$$a_n + \sum_{j=1}^k a_{n-j}((j-1)^\ell + 2\ell(j)^{\ell-1}) = (-1)^\ell - 1 \quad (3.5)$$

Proof. The local truncation error of LF-ghoRA is

$$\begin{aligned} \Delta t \tau_n = & u(t_{n+1}) - a_{n+1}u_n - (1 + a_n)u_{n-1} - \sum_{j=1}^k a_{n-j}u_{n+1-j} \\ & - 2i\omega\Delta t \left(u_n - \sum_{j=1}^k a_{n-j}u_{n-j} \right). \end{aligned} \quad (3.6)$$

Apply Taylor expansion at time t_n for any $j \geq 1$ and assume that all numerical solution at previous time are exact i.e. $u_m = u(t_m)$ for $m = 1, \dots, n$ where $t_m = m\Delta t$, then

$$u(t_{n-j}) = \sum_{\ell=0}^p \frac{(-j\Delta t)^\ell}{\ell!} u^{(\ell)}(t_n) + \mathcal{O}(\Delta t^{p+1}).$$

Substitute $u(t_{n-j})$ and $u(t_{n+1-j})$ in (3.6) and eliminate higher order terms with $\mathcal{O}(\Delta t^{p+1})$,

$$\begin{aligned}\Delta t \tau_n &= \left(- \sum_{j=-1}^k a_{n-j} \right) u_n + \left(\sum_{j=0}^k a_{n-j}(j+1) \right) \Delta t u'_n \\ &+ \sum_{\ell=2}^p \left(1 - (1+a_n)(-1)^\ell - \sum_{j=1}^k a_{n-j}((1-j)^\ell - 2\ell(-j)^{\ell-1}) \right) \frac{\Delta t^\ell}{\ell!} u_n^{(\ell)} + \mathcal{O}(\Delta t^{p+1}).\end{aligned}$$

The consistency condition (3.4) follows by setting the coefficients of u_n and u'_n equal to 0.

The additional condition (3.5) for p^{th} -order accuracy is obtained by setting the coefficient of $u_n^{(\ell)}$ to 0 for all $\ell = 2, \dots, p$, i.e.,

$$\begin{aligned}1 - (1+a_n)(-1)^\ell - \sum_{j=1}^k a_{n-j}((1-j)^\ell - 2\ell(-j)^{\ell-1}) &= 0 \\ \implies a_n + \sum_{j=1}^k a_{n-j}((j-1)^\ell + 2\ell(j)^{\ell-1}) &= (-1)^\ell - 1.\end{aligned}$$

□

3.2 ERROR ANALYSIS

In this section, we derive the phase and amplitude errors of the LF-ghoRA scheme (3.1) using the notion of modified equation, an idea related to backward error analysis. The concept is to view the numerical solution of LF-ghoRA not as an approximate solution to oscillation equation (2.1), but as an exact solution to a nearby equation, called modified equation.

Proposition 3.2.1. *The two term modified equation of oscillation equation (2.1) for p^{th} -order accurate LF-ghoRA with initial value $u(0) = x(0) = 1$ is*

$$x'(t) = i\omega x(t) + C_1(i\omega\Delta t)^p i\omega x(t) + C_2(i\omega\Delta t)^{p+1} i\omega x(t), \quad (3.7)$$

where C_1 and C_2 are defined as;

$$C_1 = \frac{\left(-1 + (1 + a_n)(-1)^{p+1} + \sum_{j=1}^k a_{n-j}((1-j)^{p+1} - 2(p+1)(-1)^p j^p) \right)}{(2 + a_n - \sum_{j=1}^k a_{n-j}(1-j))(p+1)!},$$

$$C_2 = \frac{-1 + (1 + a_n)(-1)^{p+2} + \sum_{j=1}^k a_{n-j}((1-j)^{p+2} - 2(p+2)(-1)^{p+1} j^{p+1})}{\left(2 + a_n - \sum_{j=1}^k a_{n-j}(1-j) \right) (p+2)!}$$

$$+ \frac{\left(a_n + \sum_{j=1}^k a_{n-j}(1-j)^2 \right) C_1}{2 \left(2 + a_n - \sum_{j=1}^k a_{n-j}(1-j) \right)}.$$

Proof. Consider general two term modified equation of oscillation equation for LF-ghoRA,

$$x'(t) = i\omega x(t) + \Delta t^p g_1(x(t)) + \Delta t^{p+1} g_2(x(t)).$$

Therefore,

$$x''(t) = -\omega^2 x(t) + i\omega \Delta t^p g_1(x(t)) + i\omega \Delta t^p g_1'(x(t))x(t) + \mathcal{O}(\Delta t^{p+1}),$$

$$x'''(t) = -i\omega^3 x(t) + \mathcal{O}(\Delta t^p).$$

Similarly,

$$x^{(q)}(t) = (i\omega)^q x(t) + \mathcal{O}(\Delta t^{p+3-q}) \text{ for all } q = 4, \dots, p+2.$$

The local truncation error (LTE) of two term modified equation is

$$\Delta t \tau_n = x(t_{n+1}) - a_{n+1} u_n - (1 + a_n) u_{n-1} - \sum_{j=1}^k a_{n-j} u_{n+1-j}$$

$$- 2i\omega \Delta t \left(u_n - \sum_{j=1}^k a_{n-j} u_{n-j} \right).$$

Assume that all previous numerical solution are exact, i.e. $u_m = x(t_m)$ for $m = 1 \dots n$, then

$$\Delta t \tau_n = \sum_{\ell=0}^{p+2} \frac{h^\ell}{\ell!} x^{(\ell)}(t_n) - a_{n+1} x(t_n) - (1 + a_n) \sum_{\ell=0}^{p+2} \frac{(-\Delta t)^\ell}{\ell!} x^{(\ell)}(t_n)$$

$$- \sum_{j=1}^k a_{n-j} \sum_{\ell=0}^{p+2} \frac{((1-j)\Delta t)^\ell}{\ell!} x^{(\ell)}(t_n) - 2i\omega \Delta t x(t_n)$$

$$+ 2i\omega \Delta t \sum_{j=1}^k a_{n-j} \sum_{\ell=0}^{p+2} \frac{(-j\Delta t)^\ell}{(\ell)!} x^{(\ell)}(t_n) + \mathcal{O}(\Delta t^{p+3}).$$

Use condition of p^{th} -order scheme given in (3.4) and (3.5), then all term are telescoping and cancel each other up to Δt^{p+1} . Thus, LTE reduced to

$$\begin{aligned}
\Delta t \tau_n = & \left(1 - (1 + a_n)(-1)^{p+1} - \sum_{j=1}^k a_{n-j}(1-j)^{p+1} \right. \\
& + 2(p+1)(-1)^p \sum_{j=1}^k j^p a_{n-j} \left. \right) \frac{\Delta t^{p+1}}{(p+1)!} (i\omega)^{p+1} x(t_n) \\
& + \left(2 + a_n - \sum_{j=1}^k a_{n-j}(1-j) \right) \Delta t^{p+1} g_1(x(t_n)) \\
& + \left(1 - (1 + a_n)(-1)^{p+2} - \sum_{j=1}^k a_{n-j}(1-j)^{p+2} \right. \\
& + 2(p+2)(-1)^{p+1} \sum_{j=1}^k j^{p+1} a_{n-j} \left. \right) \frac{\Delta t^{p+2}}{(p+2)!} (i\omega)^{p+2} x(t_n) \\
& + \left(-4 \sum_{j=1}^k j a_{n-j} - a_n - \sum_{j=1}^k a_{n-j}(1-j)^2 \right) i\omega \frac{\Delta t^{p+2}}{2} g_1(x(t_n)) \\
& + \left(2 + a_n - \sum_{j=1}^k a_{n-j}(1-j) \right) \Delta t^{p+2} g_2(x(t_n)) \\
& + \left(-a_n - \sum_{j=1}^k a_{n-j}(1-j)^2 \right) i\omega g_1'(x(t_n)) x(t_n) \frac{\Delta t^{p+2}}{2} + \mathcal{O}(\Delta t^{p+3}).
\end{aligned}$$

The claim follows by setting coefficient of Δt^{p+1} and Δt^{p+2} to zero. \square

Theorem 3.2.1. *The phase and amplitude error of p^{th} -order accurate LF-ghoRA is*

$$\begin{aligned}
R_+ - 1 &= \begin{cases} C_2(i\omega \Delta t)^{p+1} + \mathcal{O}(\Delta t^{p+3}), & \text{when } p \text{ is odd} \\ C_1(i\omega \Delta t)^p + \mathcal{O}(\Delta t^{p+2}), & \text{when } p \text{ is even} \end{cases} \\
|A_+| - 1 &= \begin{cases} C_1(i\omega \Delta t)^{p+1} + \mathcal{O}(\Delta t^{p+3}), & \text{when } p \text{ is odd} \\ C_2(i\omega \Delta t)^{p+2} + \mathcal{O}(\Delta t^{p+4}), & \text{when } p \text{ is even} \end{cases}
\end{aligned}$$

Proof. Consider the exact solution of oscillation and modified equation with initial value $u(0) = x(0) = 1$,

$$\begin{aligned}
u(t) &= \exp(i\omega t), \\
x(t) &= \exp(i\omega t + C_1(i\omega)^{p+1} \Delta t^p t + C_2(i\omega)^{p+2} \Delta t^{p+1} t).
\end{aligned}$$

Case 1: p is even.

The phase and amplitude error are

$$\begin{aligned} R_+ - 1 &= \frac{\arg(x(t))}{\arg(u(t))} - 1 = \frac{\omega t + C_1(i\omega)^{p+1}\Delta t^p t}{\omega t} - 1 = C_1(i\omega\Delta t)^p + \mathcal{O}(\Delta t^{p+2}) \\ |A_+| - 1 &= |x(t)| - |u(t)| = |\exp(i\omega t + C_1(i\omega)^{p+1}\Delta t^p t + C_2(i\omega)^{p+2}\Delta t^{p+1}t)| - 1 \\ &= |\exp(C_2(i\omega)^{p+2}\Delta t^{p+1}t)| - 1. \end{aligned}$$

Using approximation $\exp(C_2(i\omega)^{p+2}\Delta t^{p+1}t) \approx 1 + C_2(i\omega)^{p+2}\Delta t^{p+1}t$ since $|i\omega\Delta t| \ll 1$ and taking $t = \Delta t$ for local truncation error, we obtain

$$|A_+| - 1 = |\exp(C_2(i\omega)^{p+2}\Delta t^{p+1}\Delta t)| - 1 = C_2(i\omega\Delta t)^{p+2} + \mathcal{O}(\Delta t^{p+4}).$$

Case 2: p is odd.

The phase and amplitude error are

$$\begin{aligned} R_+ - 1 &= \frac{\arg(x(t))}{\arg(u(t))} - 1 = \frac{\omega t + C_2(i\omega)^{p+2}\Delta t^{p+1}t}{\omega t} - 1 = C_2(i\omega\Delta t)^{p+1} + \mathcal{O}(\Delta t^{p+3}) \\ |A_+| - 1 &= |x(t)| - |u(t)| = |\exp(i\omega t + C_1(i\omega)^{p+1}\Delta t^p t + C_2(i\omega)^{p+2}\Delta t^{p+1}t)| - 1 \\ &= |\exp(C_1(i\omega)^{p+1}\Delta t^p t)| - 1. \end{aligned}$$

Similarly, $\exp(C_1(i\omega)^{p+1}\Delta t^p t) \approx 1 + C_1(i\omega)^{p+1}\Delta t^p t$ since $|i\omega\Delta t| \ll 1$ and taking $t = \Delta t$ for local truncation error, we obtain

$$|A_+| - 1 = |\exp(C_1(i\omega)^{p+1}\Delta t^p\Delta t)| - 1 = C_1(i\omega\Delta t)^{p+1} + \mathcal{O}(\Delta t^{p+3}).$$

□

Corollary 3.2.1. *In particular Proposition (3.2.1) and Theorem (3.2.1) holds for second-order hoRA and fourth-order hoRA (see [32] for detail).*

Proof. First we consider the general form of second-order hoRA with $p = 2$ and $k = 2$,

$$v_{n+1} = u_{n-1} + 2\Delta t f(v_n),$$

$$u_n = v_n + a_{n+1}v_{n+1} + a_nv_n + a_{n-1}u_{n-1} + a_{n-2}u_{n-2}$$

where

$$a_{n+1} = \frac{\beta}{2}, \quad a_n = -\frac{3\beta}{2}, \quad a_{n-1} = \frac{3\beta}{2}, \quad a_{n-2} = -\frac{\beta}{2}.$$

Therefore,

$$\begin{aligned} C_1 &= \frac{\left(-1 + (1 + a_n)(-1)^{2+1} + \sum_{j=1}^2 a_{n-j}((1-j)^{2+1} - 2(2+1)(-1)^2 j^2) \right)}{(2+1)!(2 + a_n - \sum_{j=1}^2 a_{n-j}(1-j))} \\ &= \frac{(-2 - a_n - 6a_{n-1} - 25a_{n-2})}{6(2 + a_n + a_{n-2})} \\ &= \frac{5\beta - 2}{12(1 - \beta)}. \end{aligned}$$

Substitute C_1 in C_2 , we obtain,

$$\begin{aligned} C_2 &= \frac{\left(-1 + (1 + a_n)(-1)^4 + a_{n-1}(-2(4)(-1)^3 1^3 + a_{n-2}((-1)^4 - 2(4)(-1)^3 2^3) \right)}{(2 + a_n - a_{n-1}(1-1) - a_{n-2}(1-2))(4)!} \\ &\quad + \frac{(a_n + a_{n-2}(1-2)^2) C_1}{2(2 + a_n - a_{n-2}(1-2))} \\ &= \frac{a_n + 8a_{n-1} + 65a_{n-2}}{24(2 + a_n + a_{n-2})} + \frac{(a_n + a_{n-2}) C_1}{2(2 + a_n + a_{n-2})} \\ &= \frac{2\beta^2 - 3\beta}{8(1 - \beta)^2}. \end{aligned}$$

Since p is even then phase and amplitude error are

$$\begin{aligned} R_+ - 1 &= C_1(i\omega\Delta t)^2 + \mathcal{O}(\Delta t^4) = \frac{2 - 5\beta}{12(1 - \beta)}(\omega\Delta t)^2 + \mathcal{O}((\omega\Delta t)^4), \\ |A_+| - 1 &= C_2(i\omega\Delta t)^4 + \mathcal{O}(\Delta t^6) = \frac{2\beta^2 - 3\beta}{8(1 - \beta)^2}(\omega\Delta t)^4 + \mathcal{O}((\omega\Delta t)^6). \end{aligned}$$

We next consider the general form of fourth-order hoRA for $k = 3$ and $p = 4$,

$$v_{n+1} = u_{n-1} + 2\Delta t f(v_n),$$

$$u_n = v_n + a_{n+1}v_{n+1} + a_nv_n + a_{n-1}u_{n-1} + a_{n-2}u_{n-2} + a_{n-3}u_{n-3}$$

with coefficients

$$a_{n+1} = \frac{15}{53}, \quad a_n = \frac{-56}{53}, \quad a_{n-1} = \frac{78}{53}, \quad a_{n-2} = \frac{-48}{53}, \quad a_{n-3} = \frac{11}{53}.$$

Therefore C_1 and C_2 are found to be,

$$C_1 = \frac{-2 - a_n - 10a_{n-1} - 161a_{n-2} - 842a_{n-3}}{120(2 + a_n + a_{n-2} + 2a_{n-3})} = -0.8208,$$

$$C_2 = \frac{a_n + 12a_{n-1} + 385a_{n-2} + 2980a_{n-3}}{720(2 + a_n + a_{n-2} + 2a_{n-3})} + \frac{a_n + a_{n-2} + 4a_{n-3}}{2(2 + a_n + a_{n-2} + 2a_{n-3})}C_1 = 1.9045.$$

Since p is even then phase and amplitude error are

$$R_+ - 1 = -0.8208(\omega\Delta t)^4 + \mathcal{O}((\omega\Delta t)^6),$$

$$|A_+| - 1 = -1.9045(\omega\Delta t)^6 + \mathcal{O}((\omega\Delta t)^8).$$

□

Remark 3.2.1. *The phase and amplitude of second-order hoRA was founded by using Taylor expansion of physical mode and fourth-order hoRA was computed by estimating limit of physical mode while $\omega\Delta t$ approaches to zero (see [32] for detail).*

3.3 SUMMARY

We derived phase and amplitude error of pre-determined order of accuracy for general high-order Robert-Asselin time filter using notion of modified equation. In particular, we recovered error analysis for second- and fourth-order hoRA.

4.0 BACKWARD EULER PLUS FILTER

In this chapter, we have investigate the effect of simple time filter used with implicit backward Euler method. The work of this chapter is based on [22]. The fully implicit or backward Euler method is one of the first method commonly implemented when extending a code for the steady state problem and often the method of last resort for complex applications. The issue can then arise of how to increase numerical accuracy in a complex, possibly legacy code without implementing a better method. We show herein that adding one line of code to backward Euler reduces discrete curvature of the solution, increases accuracy from first to second-order, gives an immediate error estimator and induces a method akin to the second-order backward differentiation formula (BDF2). The effect of each step in the combination of 2-step is conceptually clear. The combination also extends easily to variable time step.

Consider the discretization of initial value problem (1.1) by the standard backward Euler (fully implicit) method followed by a simple time filter (next for constant time step)

$$\begin{aligned} \text{Step 1} : \quad & \frac{v_{n+1} - u_n}{\Delta t} = f(v_{n+1}), \\ \text{Step 2} : \quad & u_{n+1} = v_{n+1} - \frac{\nu}{2} \{v_{n+1} - 2u_n + u_{n-1}\}. \end{aligned} \tag{4.1}$$

where v and u denote unfiltered and once filtered values, respectively. Denote the n^{th} variable time step by Δt_n and algorithm parameter by ν . Define $t_{n+1} = t_n + \Delta t_n$ and $\tau = \Delta t_n / \Delta t_{n-1}$.

Step 2 is the only 3-point filter for which the combination of backward Euler plus a time filter produces a consistent approximation and the combination achieves second-order accuracy for $\nu = 2/3$ (Proposition 4.1.1). Tests in Section 4.4 confirm the theoretical prediction of good accuracy with $\nu = 2/3$. They also show (surprisingly) that, in this specific combination, filtering every time step reduces numerical dissipation.

The combination (4.1) is 0-stable for $-2 \leq \nu < 2$ and A -stable for $-2/3 \leq \nu \leq 2/3$ (Proposition 4.1.2). The local truncation error (LTE) for the combination when $\nu = 2/3$ is

$$LTE = -\frac{5}{6}\Delta t^3 u'''(t_n) + \mathcal{O}(\Delta t^4).$$

The constant $5/6$, while larger than the smallest error constant $1/12$ obtained from Crank-Nicolson for second-order and A -stable methods [9], is moderate. Since Step 2 with $\nu = 2/3$ has greater accuracy than Step 1, the pre- and post-filter difference

$$EST = \|u_{n+1} - v_{n+1}\| \quad (4.2)$$

can be used in a standard way as an estimator. The variable time step case is considered in Section 4.2 based on a definition of discrete curvature and the curvature reducing filter in Step 2:

$$\begin{aligned} \text{Step 1 : } & \frac{v_{n+1} - u_n}{\Delta t_n} = f(v_{n+1}), \\ \text{Step 2 : } & u_{n+1} = v_{n+1} - \frac{\nu}{2} \left\{ \frac{2\Delta t_{n-1}}{\Delta t_n + \Delta t_{n-1}} v_{n+1} - 2u_n + \frac{2\Delta t_n}{\Delta t_n + \Delta t_{n-1}} u_{n-1} \right\}. \end{aligned} \quad (4.3)$$

For second-order accuracy of (4.3), the choice of ν depends on $\tau = \Delta t_n / \Delta t_{n-1}$ and given by $\nu = \tau(1 + \tau)/(1 + 2\tau)$ (Proposition 4.2.2). The filter step reduces the discrete curvature (Definition 4.2.1) at the three points (t_{n+1}, u_{n+1}) , (t_n, u_n) , (t_{n-1}, u_{n-1}) provided $0 < \nu < 1 + \tau$ (Proposition 4.2.1). The special value of $\nu = 2/3$ for constant time step induces a one-leg, two-step, second-order accurate and strongly A -stable method which is given by

$$\frac{3}{2}u_{n+1} - 2u_n + \frac{1}{2}u_{n-1} = \Delta t f\left(\frac{3}{2}u_{n+1} - u_n + \frac{1}{2}u_{n-1}\right). \quad (4.4)$$

For general ν and variable time step, the equivalent linear multistep method is given in (4.13). The left hand side of (4.4) is the same as BDF2. The right hand side differs from BDF2 by

$$\frac{3}{2}u(t_{n+1}) - u(t_n) + \frac{1}{2}u(t_{n-1}) = u(t_{n+1}) + \mathcal{O}(\Delta t^2).$$

Remark 4.0.1. *The filter with $\nu = 1 + \tau$ is inconsistent since it forces u_{n+1} to be the linear extrapolation of u_n, u_{n-1} . Thus we assume $\nu \neq 1 + \tau$. Filtering twice is equivalent to increasing the value of the filter parameter $\nu \rightarrow \nu(2 - \frac{\nu}{2})$ and filtering once.*

Remark 4.0.2. *The two-step leapfrog scheme is commonly used in geophysical fluid dynamics simulations. The main issue it faces is an undamped computational mode, known as a time-splitting instability, [15]. Time filters centered at t_n rather than t_{n+1} , are often used in those simulations to reduce oscillations caused by this computational mode (see e.g. [2, 42, 48, 49]). The first time filter, proposed by Robert [42] and analyzed by Asselin [2], is called Robert-Asselin (RA) filter. The combination of RA filter with leapfrog is given by*

$$\text{Step 1 : } \frac{v_{n+1} - u_{n-1}}{2\Delta t} = f(v_n) \quad (\text{LF})$$

$$\text{Step 2 : } u_n = v_n + \frac{\nu}{2} \{v_{n+1} - 2v_n + u_{n-1}\}, \quad \nu \simeq 0.1. \quad (\text{RA})$$

The extension of the RA filter to variable time steps based on Section 4.2 is

$$u_n = v_n + \frac{\nu}{2} \left\{ \frac{2\Delta t_{n-1}}{\Delta t_n + \Delta t_{n-1}} v_{n+1} - 2v_n + \frac{2\Delta t_n}{\Delta t_n + \Delta t_{n-1}} u_{n-1} \right\}.$$

While this combination of leapfrog plus RA filter controls the computational mode it also overdamps the physical mode and degrades the amplitude error to first-order, [50]. Williams [48] developed a second filter step that increased the amplitude accuracy to third-order, conserved three-time-level solution mean and increased the predictability horizon of the discretized system. This combination, known as the RAW filter, is now nearly universally used in atmosphere simulations. For a one-step method, filters centered at t_n , like the RA filter, post process the computed solution but do not alter the evolution of the approximate solution. (The value of y_n is changed after u_n is used to calculate u_{n+1} . The changed value is not used thereafter.) For that reason the filter is shifted to t_{n+1} herein for use with single step methods.

4.1 CONSTANT TIME STEP

We develop the properties of the method for constant time step in this section.

4.1.1 Derivation of the Method

Consider backward Euler plus a general 3-point time filter

$$\text{Step 1 : } \frac{v_{n+1} - u_n}{\Delta t} = f(v_{n+1}), \quad (4.5)$$

$$\text{Step 2 : } u_{n+1} = v_{n+1} + \{av_{n+1} + bu_n + cu_{n-1}\}.$$

Proposition 4.1.1. *Let the time step be constant and $\nu \neq 2$; (4.5) is consistent if and only if $a = -\frac{\nu}{2}$, $c = -\frac{\nu}{2}$, $b = \nu$ for some ν . Thus, the step 2 of (4.5) is*

$$u_{n+1} = v_{n+1} - \frac{\nu}{2} (v_{n+1} - 2u_n + u_{n-1}). \quad (4.6)$$

In this case, the combination of 2-step gives following equivalent linear multistep method;

$$\begin{aligned} \frac{1}{1 - \frac{\nu}{2}} u_{n+1} - \frac{1 + \frac{\nu}{2}}{1 - \frac{\nu}{2}} u_n + \frac{\frac{\nu}{2}}{1 - \frac{\nu}{2}} u_{n-1} &= \\ &= \Delta t f\left(\frac{1}{1 - \frac{\nu}{2}} u_{n+1} - \frac{\nu}{1 - \frac{\nu}{2}} u_n + \frac{\frac{\nu}{2}}{1 - \frac{\nu}{2}} u_{n-1}\right). \end{aligned} \quad (4.7)$$

The combination is second-order accurate if and only if $\nu = 2/3$.

Proof. Eliminating $v_{n+1} = \frac{1}{1+a} (u_{n+1} - bu_n - cu_{n-1})$ yields the equivalent one-leg, linear multistep method for the post-filtered values

$$\frac{1}{1+a} u_{n+1} - \frac{1+a+b}{1+a} u_n - \frac{c}{1+a} u_{n-1} = \Delta t f\left(\frac{1}{1+a} u_{n+1} - \frac{b}{1+a} u_n - \frac{c}{1+a} u_{n-1}\right).$$

A Taylor series calculation shows that consistency requires

$$a + b + c = 0 \text{ and } a = c.$$

Thus, the method is consistent when

$$c = a \text{ and } b = -2a.$$

Taking $a = \nu/2$ gives (4.6) and (4.7). The same Taylor series expansion for second-order accuracy yields $c = -1/3$ corresponding to $\nu = 2/3$ as claimed. \square

4.1.2 Stability for Constant Time Step

The general two-step method obtained from (1.2) with $k = 2$ is given by

$$\alpha_0 u_{n+1} + \alpha_1 u_n + \alpha_2 u_{n-1} = \Delta t f(\beta_0 u_{n+1} + \beta_1 u_n + \beta_2 u_{n-1}). \quad (4.8)$$

The equivalent linear multistep method (4.7) corresponds to (4.8) with coefficients

$$\begin{aligned} \alpha_0 &= \frac{1}{1 - \frac{\nu}{2}}, \quad \alpha_1 = -\frac{1 + \frac{\nu}{2}}{1 - \frac{\nu}{2}}, \quad \alpha_2 = \frac{\frac{\nu}{2}}{1 - \frac{\nu}{2}}, \\ \beta_0 &= \frac{1}{1 - \frac{\nu}{2}}, \quad \beta_1 = -\frac{\nu}{1 - \frac{\nu}{2}}, \quad \beta_2 = \frac{\frac{\nu}{2}}{1 - \frac{\nu}{2}}. \end{aligned}$$

The A -stability of general two-step methods (4.8) are characterized in terms of their coefficients in [10, 11, 13, 21, 37]. In this section and the following one we apply the characterization for variable time step in [11], which states that the method is A -stable if

$$\begin{cases} -\alpha_1 \geq 0, \\ 1 - 2\beta_1 \geq 0, \\ 2(\beta_0 - \beta_2) + \alpha_1 \geq 0. \end{cases} \quad (4.9)$$

Proposition 4.1.2. *The linear multistep method (4.7) is 0-stable for $-2 \leq \nu < 2$ and A -stable for*

$$-\frac{2}{3} \leq \nu \leq \frac{2}{3}.$$

Proof. For 0-stability, the associated polynomial

$$\frac{1}{1 - (\nu/2)} z^2 - \frac{1 + (\nu/2)}{1 - (\nu/2)} z + \frac{(\nu/2)}{1 - (\nu/2)} = 0$$

has roots

$$z_{\pm} = 1, \quad \frac{\nu}{2}$$

from which 0-stability follows for $-2 \leq \nu < 2$.

We apply the characterization (4.9) to show A -stability for $-2 \leq \nu < 2$. A short calculation show the first condition $-\alpha_1 \geq 0$ holds if and only if $\nu \geq -2$. The second condition $1 - 2\beta_1 \geq 0$ holds if and only if $\nu \geq -\frac{2}{3}$. The third condition $2(\beta_0 - \beta_2) + \alpha_1 \geq 0$ holds if and only if $\nu \leq \frac{2}{3}$. \square

For $\nu = 2/3$, the stability region of (4.1) is computed by root locus curve method, e.g., [19, 24], and presented in Figure 4.1. For comparison, Figure 4.2 presents the stability region boundaries of backward Euler, backward Euler plus filter and BDF2. The boundaries of the stability regions in Figure 4.2 suggest that all three methods comparably dissipative with backward Euler plus filter the most dissipative. This is incorrect. Tests with periodic and quasi-periodic solution in Section 4.4 show that BDF2 and backward Euler plus filter are comparably dissipative. In these tests, both are significantly less dissipative than backward Euler.

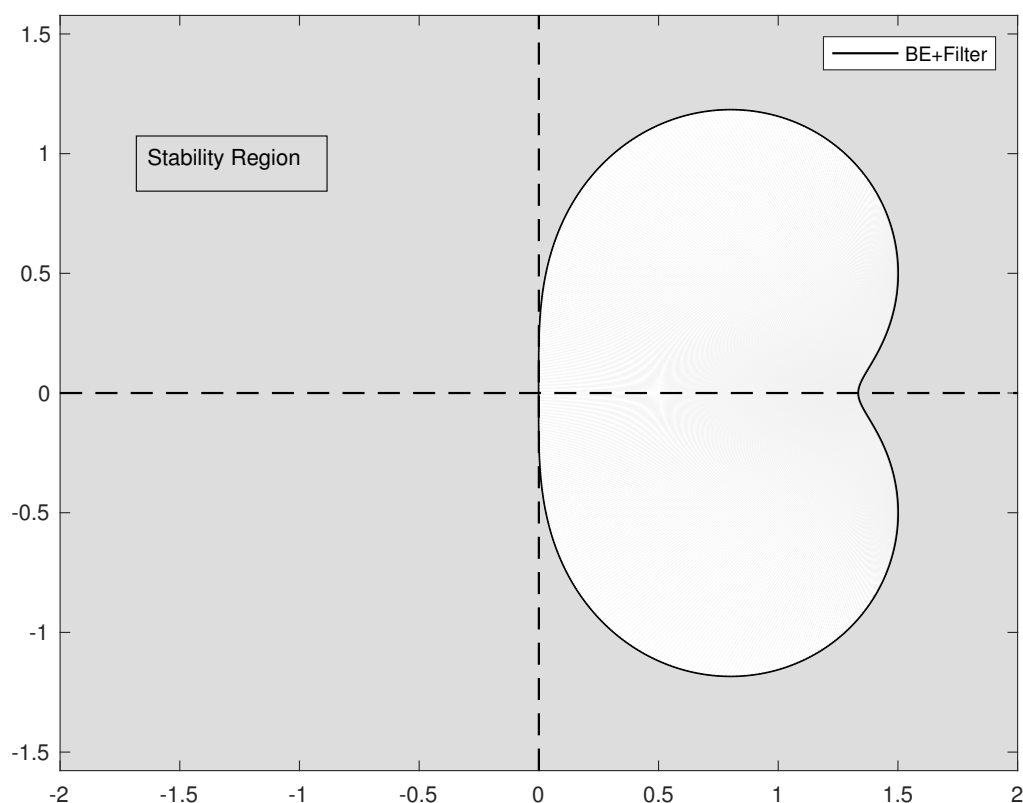


Figure 4.1: Stability region of backward Euler plus time filter

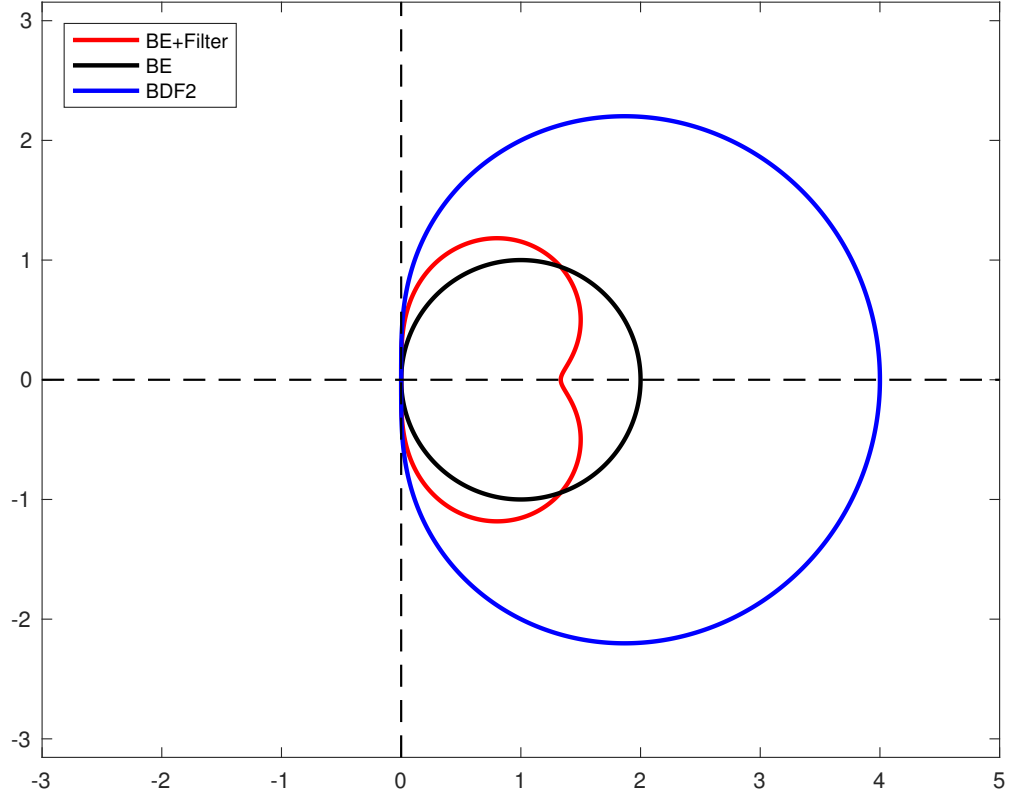


Figure 4.2: Boundaries of Stability Regions

4.2 VARIABLE TIME STEP

Since backward Euler is a one-step method its extension to variable time step is clear. Thus, the key to a variable time step realization of the method (4.5) is to extend the time filter (4.6) to variable time step.

4.2.1 Time Filter for Variable Time Step

To extend time filters to variable time step, we must first define the discrete curvature. The extension of differential geometry to discrete setting is an active research field with

considerable work on discrete curvature, e.g., [36]. For three points the natural definitions are either the discrete second difference or the inverse of the radius of the interpolating circle. We employ the former scaled by $\Delta t_{n-1}\Delta t_n$ in order to be consistent with work in geophysical fluid dynamics [29, 49]. Let the standard Lagrange basis functions for t_{n-1} , t_n , t_{n+1} be $\ell_{n-1}(t)$, $\ell_n(t)$, $\ell_{n+1}(t)$, respectively. The quadratic interpolant polynomial ϕ , from which a discrete curvature is calculated, through the points (t_{n-1}, u_{n-1}) , (t_n, u_n) , (t_{n+1}, u_{n+1}) is

$$\phi(t) = u_{n+1}\ell_{n+1}(t) + u_n\ell_n(t) + u_{n-1}\ell_{n-1}(t).$$

Definition 4.2.1. *The discrete curvature at (t_{n-1}, u_{n-1}) , (t_n, u_n) , and (t_{n+1}, u_{n+1}) is*

$$\begin{aligned}\kappa &= \Delta t_{n-1}\Delta t_n\phi'' = \frac{2\Delta t_{n-1}}{\Delta t_n + \Delta t_{n-1}}u_{n+1} - 2u_n + \frac{2\Delta t_n}{\Delta t_n + \Delta t_{n-1}}u_{n-1} \\ &= \frac{2}{1+\tau}u_{n+1} - 2u_n + \frac{2\tau}{1+\tau}u_{n-1}.\end{aligned}$$

The extension of (4.6) to variable time step is

$$u_{n+1} = v_{n+1} - \frac{\nu}{2} \left\{ \frac{2}{1+\tau}v_{n+1} - 2u_n + \frac{2\tau}{1+\tau}u_{n-1} \right\}. \quad (4.10)$$

Proposition 4.2.1 shows that the step (4.10) is curvature reducing. This establishes that Step 2 does not introduce oscillatory instabilities.

Proposition 4.2.1. *For (4.10) the discrete curvature before, κ^{old} , and after, κ^{new} , filtering satisfies*

$$\kappa^{new} = \left(1 - \frac{\nu}{1+\tau}\right)\kappa^{old}.$$

It reduces, without changing sign, the discrete curvature, $|\kappa^{new}| < |\kappa^{old}|$, provided

$$0 < \nu < 1 + \tau.$$

Proof. First multiply (4.10) through by $\frac{2}{1+\tau}$, we get

$$\frac{2}{1+\tau}u_{n+1} = \frac{2}{1+\tau}v_{n+1} - \frac{\nu}{2} \cdot \frac{2}{1+\tau} \left\{ \frac{2}{1+\tau}v_{n+1} - 2u_n + \frac{2\tau}{1+\tau}u_{n-1} \right\}. \quad (4.11)$$

Then, adding $-2u_n + \frac{2\tau}{1+\tau}u_{n-1}$ to both sides of (4.11) gives

$$\frac{2}{1+\tau}u_{n+1} - 2u_n + \frac{2\tau}{1+\tau}u_{n-1} = \left(1 - \frac{\nu}{1+\tau}\right) \left\{ \frac{2}{1+\tau}v_{n+1} - 2u_n + \frac{2\tau}{1+\tau}u_{n-1} \right\}.$$

Thus, we obtain

$$\kappa^{new} = \left(1 - \frac{\nu}{1+\tau}\right) \kappa^{old}.$$

Curvature reduction holds provided $0 < \frac{\nu}{1+\tau} < 1$, as claimed. \square

4.2.2 The Local Truncation Error

Since Step 1 and Step 2 are now well defined for variable time step, the method is determined.

It is as follows:

$$\begin{aligned} \text{Step 1 : } \quad & \frac{v_{n+1} - u_n}{\Delta t_n} = f(v_{n+1}), \\ \text{Step 2 : } \quad & u_{n+1} = v_{n+1} - \frac{\nu}{2} \left\{ \frac{2}{1+\tau}v_{n+1} - 2u_n + \frac{2\tau}{1+\tau}u_{n-1} \right\}. \end{aligned} \quad (4.12)$$

Step 2 in (4.12) is used to solve for $v_{n+1} = \frac{1+\tau}{1+\tau-\nu}u_{n+1} - \nu\frac{1+\tau}{1+\tau-\nu}u_n + \frac{\tau\nu}{1+\tau-\nu}u_{n-1}$ and eliminate v_{n+1} in Step 1. This gives the equivalent two-step method

$$\begin{aligned} \frac{1+\tau}{1+\tau-\nu}u_{n+1} - \nu\frac{1+\tau}{1+\tau-\nu}u_n + \frac{\tau\nu}{1+\tau-\nu}u_{n-1} - u_n = \\ = \Delta t_n f\left(\frac{1+\tau}{1+\tau-\nu}u_{n+1} - \nu\frac{1+\tau}{1+\tau-\nu}u_n + \frac{\tau\nu}{1+\tau-\nu}u_{n-1}\right). \end{aligned} \quad (4.13)$$

This corresponds to a linear multistep method (4.8) with coefficients

$$\begin{aligned} \alpha_0 &= \frac{1+\tau}{1+\tau-\nu}, \quad \alpha_1 = -\frac{1+\tau+\nu\tau}{1+\tau-\nu}, \quad \alpha_2 = \frac{\tau\nu}{1+\tau-\nu}, \\ \beta_0 &= \frac{1+\tau}{1+\tau-\nu}, \quad \beta_1 = -\nu\frac{1+\tau}{1+\tau-\nu}, \quad \beta_2 = \frac{\tau\nu}{1+\tau-\nu}. \end{aligned} \quad (4.14)$$

Proposition 4.2.2. *The variable time step method (4.13) is consistent for $\nu \neq 1 + \tau$. It is second-order accurate provided*

$$\nu = \frac{\tau(1 + \tau)}{1 + 2\tau}. \quad (4.15)$$

Moreover, the local truncation error (LTE) for $\nu = \frac{\tau(\tau+1)}{1+2\tau}$ is

$$LTE = -\frac{(1 + 4\tau)}{6\tau} \Delta t_n^3 u'''(t_n) + \mathcal{O}(\Delta t_n^4). \quad (4.16)$$

Proof. By a Taylor expansion, the method is consistent if and only if

$$\frac{1 + \tau}{1 + \tau - \nu} + \left(-\nu \frac{1 + \tau}{1 + \tau - \nu} - 1\right) + \frac{\tau\nu}{1 + \tau - \nu} = 0,$$

and

$$\frac{1 + \tau}{1 + \tau - \nu} - \frac{1}{\tau} \frac{\tau\nu}{1 + \tau - \nu} - \left(\frac{1 + \tau}{1 + \tau - \nu} - \nu \frac{1 + \tau}{1 + \tau - \nu} + \frac{\tau\nu}{1 + \tau - \nu}\right) = 0.$$

These two consistency conditions identically hold. By the same expansion, the method is second-order accurate if and only if

$$\tau^2 \frac{1 + \tau}{1 + \tau - \nu} + \frac{\tau\nu}{1 + \tau - \nu} - 2\tau^2 \frac{1 + \tau}{1 + \tau - \nu} + 2\tau \frac{\tau\nu}{1 + \tau - \nu} = 0.$$

This holds if and only if ν is given by

$$\nu = \frac{\tau + \tau^2}{1 + 2\tau}. \quad (4.17)$$

For $\nu = \frac{\tau(\tau+1)}{1+2\tau}$ the leading term in the LTE is calculated to be (4.16). □

Remark 4.2.1. *The usual variable time step BDF2 method, [3, 17], is given by*

$$\frac{2\tau + 1}{\tau + 1} u_{n+1} - (\tau + 1) u_n + \frac{\tau^2}{\tau + 1} u_{n-1} = \Delta t_n f(u_{n+1}). \quad (4.18)$$

For $\nu = \tau(1 + \tau)/(1 + 2\tau)$ the left hand side of (4.13) is the same as variable time step BDF2 while the right hand side again differs.

4.2.3 Stability for Variable Time Step

As defined by Dahlquist, Liniger and Nevanlinna [13] in equation (1.12) p.1072, a variable time step method is A -stable if, when applied as a one-leg scheme to

$$u' = \lambda(t)u, \text{ where } \operatorname{Re}(\lambda(t)) \leq 0,$$

solutions are always bounded for any sequence of time step and any such $\lambda(t)$. Conditions (4.9) were derived in [11] for variable time step. We therefore analyze A -stability for variable time step applying conditions (4.9).

Proposition 4.2.3. *The method (4.13) is A -stable for*

$$-\frac{1+\tau}{1+2\tau} \leq \nu \leq \min\left\{\frac{1+\tau}{3\tau}, 1+\tau\right\}.$$

Proof. We check conditions (4.9) for the coefficients (4.14). The first condition holds if

$$\frac{1+\tau+\nu\tau}{1+\tau-\nu} \geq 0.$$

This holds if and only if

$$-\frac{1+\tau}{\tau} \leq \nu \leq 1+\tau.$$

The second condition holds if

$$1+2\nu\frac{1+\tau}{1+\tau-\nu} \geq 0.$$

This holds if and only if

$$-\frac{1+\tau}{1+2\tau} \leq \nu.$$

As $\nu \leq 1+\tau$ condition 3 holds if

$$2\left(\frac{1+\tau}{1+\tau-\nu} - \frac{\tau\nu}{1+\tau-\nu}\right) - \frac{1+\tau+\nu\tau}{1+\tau-\nu} \geq 0.$$

This condition holds if and only if

$$\nu \leq \frac{1+\tau}{3\tau}.$$

Since $\min\left\{\frac{1+\tau}{\tau}, \frac{1+\tau}{1+2\tau}\right\} = \frac{1+\tau}{1+2\tau}$ the result follows. □

Since the filter is curvature reducing only for $0 < \nu < 1 + \tau$, it is sensible to restrict the values to $0 < \nu \leq \min\{\frac{1+\tau}{3\tau}, 1 + \tau\}$. We plot next this region in Figure 4.3. The region of variable step A -stability is below the dark curve in the (τ, ν) plane. The dashed curve plotted is the choice of $\nu = \nu(\tau)$ that yields second-order accuracy. The method is A -stable for constant or decreasing time step. If increasing the time step, one must either accept first-order accuracy with A -stability or second-order with some reduced (and yet undetermined) $A(\alpha)$ -stability for $\alpha < \pi/2$.

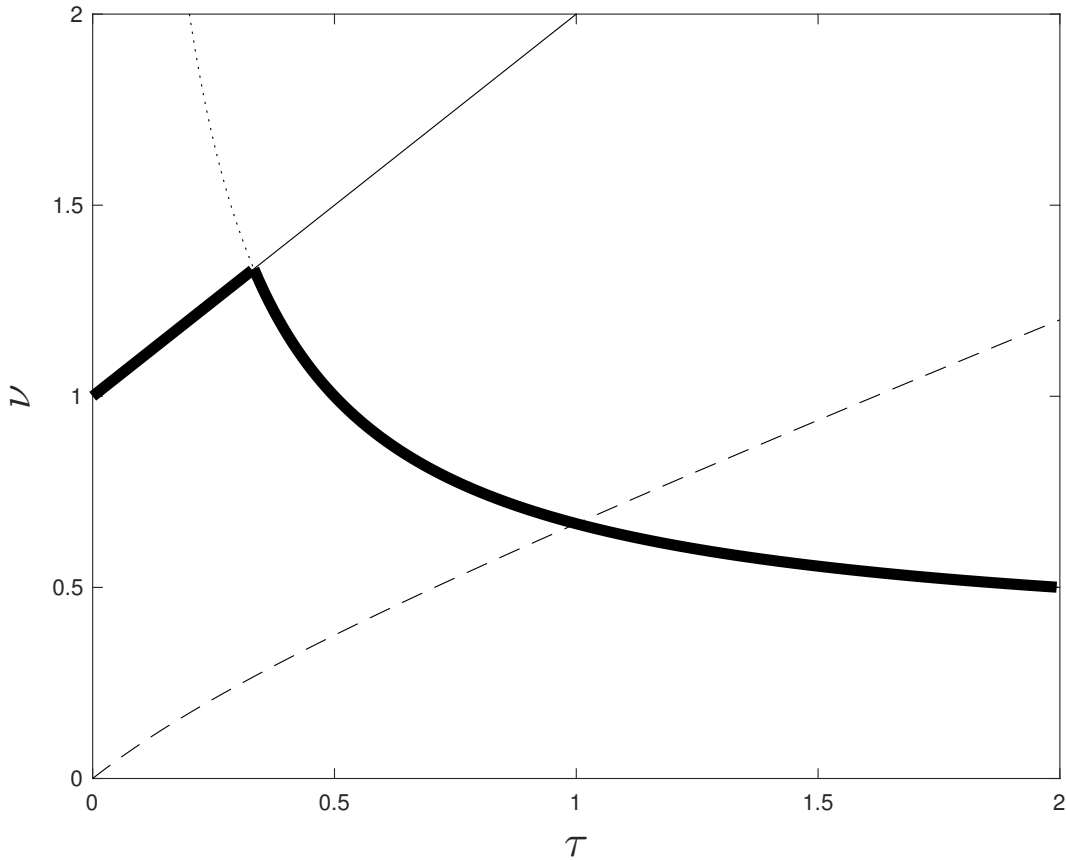


Figure 4.3: A -stable for $\nu \leq$ dark curve, Dashed Curve = $\mathcal{O}(\Delta t^2)$

Remark 4.2.2. *Variable time step BDF2 is A -stable for $\tau \leq 1$. This is the same constraint as for the method herein when ν chosen for second-order accuracy.*

4.2.4 Adaptive Time Step Algorithm

The combination of backward Euler plus filter lends itself to adaptive implementation. There are various choices that must be made in such an implementation(see [6, 18, 39, 43]). We have purposefully made the simplest one of each option. The general adaptive method with simple halving and doubling time step, safety factor s and order p implemented as follows¹.

Algorithm 4.1: Halving and Doubling Time step

```

Given      :  $u_n, u_{n-1}, \Delta t_n, \Delta t_{n-1}, Tol$  and  $tFinal$ ;

while  $t_n < tFinal$  do
    •  $t_{n+1} = t_n + \Delta t_n$ ;
    •  $\nu = \frac{\Delta t_n(\Delta t_n + \Delta t_{n-1})}{\Delta t_{n-1}(2\Delta t_n + \Delta t_{n-1})}$ ;
    •  $v_{n+1} - u_n = \Delta t_n f(v_{n+1})$ ;
    •  $u_{n+1} = v_{n+1} - \frac{\nu}{2} \left\{ \frac{2\Delta t_{n-1}}{\Delta t_n + \Delta t_{n-1}} v_{n+1} - 2u_n + \frac{2\Delta t_n}{\Delta t_n + \Delta t_{n-1}} u_{n-1} \right\}$ ;
    •  $Est = \|u_{n+1} - v_{n+1}\|$ ;

    if  $Tol < s * Est$  then
        |  $\Delta t_n = \frac{\Delta t_n}{2}$ ;
        |  $n = n$ ;
    else if  $Est \leq s * \frac{Tol}{2^{p+1}}$  then
        |  $\Delta t_{n+1} = 2\Delta t_n$ ;
        |  $n = n + 1$ ;
    else
        |  $\Delta t_{n+1} = \Delta t_n$ ;
        |  $n = n + 1$ ;
    end
end

```

¹The safety factor s is usually taken around 0.95.

4.3 ERROR ANALYSIS FOR PHASE AND AMPLITUDE

In this section, we use modified equation to derive phase and amplitude error of backward Euler plus filter. The linear multistep method (4.13) is generally a first-order approximation to oscillation equation (2.1) and second-order for the choice $\nu = \tau(1 + \tau)/(1 + 2\tau)$. To delineate the distribution of error between phase error and amplitude error we construct the modified equation of the method for the oscillation equation. We note that the modified equation is based on an expansion that assumes implicit condition $|\omega\Delta t_n| \ll 1$.

Proposition 4.3.1. *The three term modified equation of oscillation equation (2.1) with initial value $u(0) = 1$ for linear multistep method (4.13) is*

$$\begin{aligned} x'(t) &= i\omega x(t) + \Delta t_n C_1 (i\omega)^2 x(t) + \Delta t_n^2 C_2 (i\omega)^3 x(t) + \Delta t_n^3 C_3 (i\omega)^4 x(t), \\ x(0) &= 1. \end{aligned} \tag{4.19}$$

where C_1, C_2, C_3 , are

$$\begin{aligned} C_1 &= \frac{\tau + \tau^2 - \nu - 2\nu\tau}{2\tau(1 + \tau - \nu)}, \\ C_2 &= \frac{2\tau^4 + \tau^3(4 - 5\nu) + \nu(1 + 2\nu) + \tau\nu(1 + 6\nu) + \tau^2(2 - 5\nu + 6\nu^2)}{6\tau^2(1 + \tau - \nu)^2}, \\ C_3 &= \frac{6\tau^6 + \tau^5(18 - 20\nu) - \tau\nu^2(31 + 24\nu) + \tau^2\nu(4 - 33\nu - 36\nu^2)}{24\tau^3(1 + \tau - \nu)^3} \\ &\quad + \frac{\tau^4(18 - 39\nu + 23\nu^2) + \tau^3(6 - 16\nu + 13\nu^2 - 24\nu^3) - \nu(1 + 8\nu + 6\nu^2)}{24\tau^3(1 + \tau - \nu)^3}. \end{aligned}$$

Proof. The general three term modified equation of oscillation equation with initial value $x(0) = 1$ takes the form

$$x' = i\omega x + \Delta t_n g_1(x) + \Delta t_n^2 g_2(x) + \Delta t_n^3 g_3(x).$$

Thus,

$$\begin{aligned} x'' &= -\omega^2 x + i\omega\Delta t_n g_1(x) + i\omega\Delta t_n^2 g_2(x) \\ &\quad + i\omega\Delta t_n g_1'(x)x + \Delta t_n^2 g_1'(x)g_1(x) + i\omega\Delta t_n^2 g_2'(x)x + \mathcal{O}(\Delta t_n^3), \\ x''' &= -i\omega^3 x - \omega^2\Delta t_n g_1(x) - 2\omega^2\Delta t_n g_1'(x)x + \mathcal{O}(\Delta t_n^2), \\ x^{(4)} &= \omega^4 x + \mathcal{O}(\Delta t_n). \end{aligned}$$

Consider linear multistep method (4.13) applied to oscillation equation (2.1),

$$u_{n+1} - \frac{\nu\tau + 1 + \tau}{\tau + 1}u_n + \frac{\nu\tau}{1 + \tau}u_{n-1} = i\omega\Delta t_n u_{n+1} - i\omega\Delta t_n \nu u_n + \frac{\nu\tau}{1 + \tau}i\omega\Delta t_n u_{n-1}.$$

Rearrange term and eliminate u_{n+1} , we obtain

$$u_{n+1} = \frac{1}{1 - i\omega\Delta t_n} \left(\frac{\nu\tau + 1 + \tau}{\tau + 1}u_n - \frac{\nu\tau}{1 + \tau}u_{n-1} - i\omega\Delta t_n \nu u_n + \frac{\nu\tau}{1 + \tau}i\omega\Delta t_n u_{n-1} \right).$$

We can use approximation of $\frac{1}{1 - i\omega\Delta t_n} \approx 1 + i\omega\Delta t_n - \omega^2\Delta t_n^2 - i\omega^3\Delta t_n^3 + \omega^4\Delta t_n^4 + \mathcal{O}(\Delta t_n^5)$ since $|i\omega\Delta t_n| \ll 1$. Therefore,

$$\begin{aligned} u_{n+1} = & \frac{\nu\tau + 1 + \tau}{\tau + 1}u_n - \frac{\nu\tau}{1 + \tau}u_{n-1} + \left(\frac{1 + \tau - \nu}{\tau + 1} \right) i\omega\Delta t_n u_n + \left(\frac{\nu - 1 - \tau}{\tau + 1} \right) \omega^2\Delta t_n^2 u_n \\ & + \left(\frac{\nu - 1 - \tau}{\tau + 1} \right) i\omega^3\Delta t_n^3 u_n + \left(\frac{1 + \tau - \nu}{\tau + 1} \right) \omega^4\Delta t_n^4 u_n + \mathcal{O}(\Delta t_n^5). \end{aligned}$$

The local truncation error of variable time step method (4.13) with modified equation is

$$\begin{aligned} LTE = & x(t_{n+1}) - u_{n+1} = x(t_{n+1}) - \frac{\nu\tau + 1 + \tau}{\tau + 1}u_n + \frac{\nu\tau}{1 + \tau}u_{n-1} \\ & - \left(\frac{1 + \tau - \nu}{\tau + 1} \right) i\omega\Delta t_n u_n - \left(\frac{\nu - 1 - \tau}{\tau + 1} \right) \omega^2\Delta t_n^2 u_n \\ & - \left(\frac{\nu - 1 - \tau}{\tau + 1} \right) i\omega^3\Delta t_n^3 u_n - \left(\frac{1 + \tau - \nu}{\tau + 1} \right) \omega^4\Delta t_n^4 u_n + \mathcal{O}(\Delta t_n^5). \end{aligned}$$

Assume that numerical solution of all previous time steps are exact, i.e., $u_i = x(t_i)$ for all $i = 1 \cdots n$,

$$\begin{aligned} LTE = & x(t_{n+1}) - \frac{\nu\tau + 1 + \tau}{\tau + 1}x(t_n) + \frac{\nu\tau}{1 + \tau}x(t_{n-1}) \\ & - \left(\frac{1 + \tau - \nu}{\tau + 1} \right) i\omega\Delta t_n x(t_n) - \left(\frac{\nu - 1 - \tau}{\tau + 1} \right) \omega^2\Delta t_n^2 x(t_n) \\ & - \left(\frac{\nu - 1 - \tau}{\tau + 1} \right) i\omega^3\Delta t_n^3 x(t_n) - \left(\frac{1 + \tau - \nu}{\tau + 1} \right) \omega^4\Delta t_n^4 x(t_n) + \mathcal{O}(\Delta t_n^5). \end{aligned}$$

Apply the Taylor expansion of $x(t_{n-1})$, $x(t_{n+1})$ at time t_n and substitute $x(t_{n+1})$, $x(t_{n-1})$, $x'(t_n)$, $x''(t_n)$, $x'''(t_n)$ and $x^{(4)}(t_n)$ in LTE , we get

$$\begin{aligned}
LTE = & \left[\left(\frac{1+\tau-\nu}{1+\tau} \right) g_1(x(t_n)) - \frac{1}{2} \omega^2 x(t_n) - \frac{\nu}{2(\tau+\tau^2)} \omega^2 x(t_n) - \left(\frac{\nu-1-\tau}{\tau+1} \right) \omega^2 x(t_n) \right] \Delta t_n^2 \\
& + \left[\left(\frac{1}{2} + \frac{\nu}{2(\tau+\tau^2)} \right) i\omega g_1(x(t_n)) + \left(\frac{1}{2} + \frac{\nu}{2(\tau+\tau^2)} \right) i\omega g_1'(x(t_n)) x(t_n) \right. \\
& - \left(\frac{1}{6} - \frac{\nu}{6(\tau^2+\tau^3)} \right) i\omega^3 x(t_n) - \left(\frac{\nu-1-\tau}{\tau+1} \right) i\omega^3 x(t_n) + \frac{1+\tau-\nu}{1+\tau} g_2(x(t_n)) \left. \right] \Delta t_n^3 \\
& + \left[\frac{1+\tau-\nu}{1+\tau} g_3(x(t_n)) + \left(\frac{1}{2} + \frac{\nu}{2(\tau+\tau^2)} \right) i\omega g_2(x(t_n)) \right. \\
& + \left(\frac{1}{2} + \frac{\nu}{2(\tau+\tau^2)} \right) g_1(x(t_n)) g_1'(x(t_n)) + \left(\frac{1}{2} + \frac{\nu}{2(\tau+\tau^2)} \right) i\omega g_2'(x(t_n)) x(t_n) \\
& - \left(\frac{1}{6} - \frac{\nu}{6(\tau^2+\tau^3)} \right) \omega^2 g_1(x(t_n)) - \left(\frac{1}{6} - \frac{\nu}{6(\tau^2+\tau^3)} \right) 2\omega^2 g_1'(x(t_n)) x(t_n) \\
& \left. + \left(\frac{1}{24} + \frac{\nu}{24(\tau^3+\tau^4)} \right) \omega^4 x(t_n) - \left(\frac{1+\tau-\nu}{\tau+1} \right) \omega^4 x(t_n) \right] \Delta t_n^4.
\end{aligned}$$

Setting coefficient of Δt_n^2 term equal to zero to find $g_1(x)$

$$g_1(x) = \frac{2\nu\tau - \tau - \tau^2 + \nu}{2\tau(1+\tau-\nu)} \omega^2 x = C_1(i\omega)^2 x.$$

We use $g_1(x)$ and set coefficient of Δt_n^3 equal to zero, we obtain $g_2(x)$ as following,

$$g_2(x) = -\frac{2\tau^4 + \tau^3(4-5\nu) + \nu(1+2\nu) + \tau\nu(1+6\nu) + \tau^2(2-5\nu+6\nu^2)}{6\tau^2(1+\tau-\nu)^2} i\omega^3 x = C_2(i\omega)^3 x.$$

Finally, we use $g_1(x)$ and $g_2(x)$ and set coefficient of Δt_n^4 to zero, we get

$$\begin{aligned}
g_3(x) = & \left[\frac{6\tau^6 + \tau^5(18-20\nu) - \tau\nu^2(31+24\nu) + \tau^2\nu(4-33\nu-36\nu^2)}{24\tau^3(1+\tau-\nu)^3} \right. \\
& + \left. \frac{\tau^4(18-39\nu+23\nu^2) + \tau^3(6-16\nu+13\nu^2-24\nu^3) - \nu(1+8\nu+6\nu^2)}{24\tau^3(1+\tau-\nu)^3} \right] \omega^4 x \\
= & C_4(i\omega)^4 x.
\end{aligned}$$

□

Remark 4.3.1. The linear multistep method (4.13) of variable time step is generally a first-order approximation oscillation equation (2.1) and fourth-order approximation to modified equation (4.19).

Theorem 4.3.1. The phase and amplitude error of linear multistep method (4.13) for variable time step is

$$R_+ - 1 = -C_2(\omega\Delta t_n)^2 + \mathcal{O}((\omega\Delta t_n)^4)$$

$$|A_+| - 1 = -C_1(\omega\Delta t_n)^2 + C_3(\omega\Delta t_n)^4 + \mathcal{O}((\omega\Delta t_n)^6).$$

where C_1, C_2, C_3 are defined in (4.3.1).

Proof. Consider exact solution of oscillation equation (2.1) and modified equation (4.19) with initial value $u(0) = x(0) = 1$,

$$u(t) = \exp(i\omega t) = \cos(\omega t) + i \sin(\omega t).$$

$$x(t) = \exp(i\omega t + \Delta t_n C_1 (i\omega)^2 t + \Delta t_n^2 C_2 (i\omega)^3 t + \Delta t_n^3 C_3 (i\omega)^4 t)$$

$$= \exp(-\Delta t_n C_1 \omega^2 t + \Delta t_n^3 C_3 \omega^4 t) [\cos(\omega t - \Delta t_n^2 C_2 \omega^3 t) + i \sin(\omega t - \Delta t_n^2 C_2 \omega^3 t)].$$

Thus, phase error is

$$R_+ - 1 = \frac{\arg(x(t))}{\arg(u(t))} - 1 = \frac{\omega t - \Delta t_n^2 C_2 \omega^3 t}{\omega t} - 1 = -C_2(\omega\Delta t_n)^2,$$

and amplitude error is

$$|A_+| - 1 = \exp(-\Delta t_n C_1 \omega^2 t + \Delta t_n^3 C_3 \omega^4 t) - 1.$$

Since we are looking for local error in one step, take $t = \Delta t_n$ and use approximation

$$\exp(-\Delta t_n^2 C_1 \omega^2 + \Delta t_n^4 C_3 \omega^4) \approx 1 - C_1(\omega\Delta t_n)^2 + C_3(\omega\Delta t_n)^4.$$

Thus,

$$|A_+| - 1 = -C_1(\omega\Delta t_n)^2 + C_3(\omega\Delta t_n)^4.$$

□

The comparison of phase speed and amplitude of backward Euler and backward euler plus filter is presented in Figure 4.4.

Remark 4.3.2. The linear multistep method (4.13) with variable time step is second-order accurate and fourth-order accurate in amplitude when $C_1 = 0$ i.e. $\nu = \frac{\tau(\tau+1)}{2\tau+1}$.

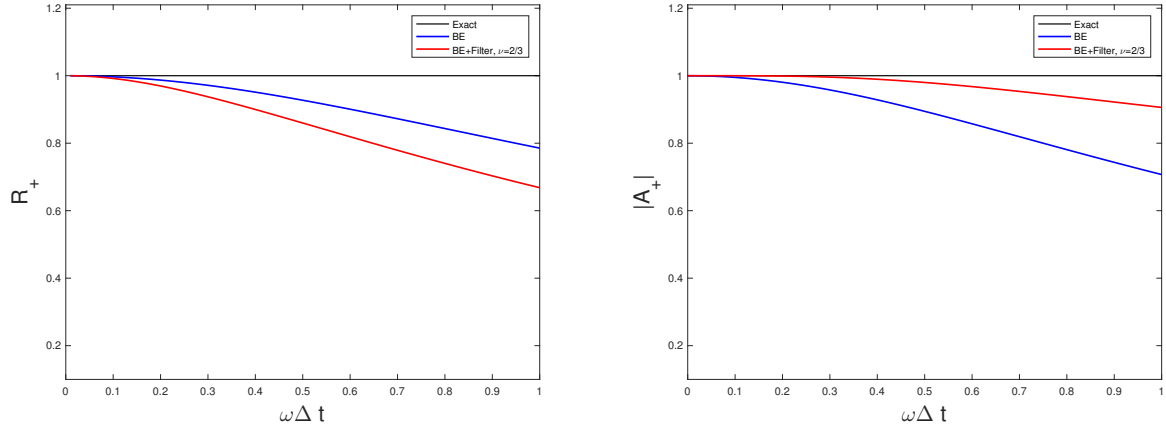


Figure 4.4: Comparison of phase speed and amplitude of physical mode for backward Euler method and second-order backward Euler plus filter method.

4.4 NUMERICAL TESTS

We present a few numerical illustrations. Where appropriate the RKF4-5 solution is used as benchmark. The first tests is for the Lorenz system. First backward Euler plus filter ($\nu = 2/3$) is compared with backward Euler (Step 1 without Step 2) and BDF2. The second test, from [44], is an example of one for which backward Euler preserves Lyapunov stability of the steady state while common variants do not. The third test is for a periodic (nonlinear pendulum) and a quasi-periodic oscillation. The fourth test is for the Van der Pol equation [45] with parameter $\mu = 1000$. This is a classic test problem for stiff solvers.

4.4.1 The Lorenz System

Consider the Lorenz system from [35]

$$\begin{aligned}\frac{dX}{dt} &= 10(Y - X), \\ \frac{dY}{dt} &= -XZ + 28X - Y, \\ \frac{dZ}{dt} &= XY - \frac{8}{3}Z.\end{aligned}$$

We use the standard parameter values of Lorenz [35]. These produce a chaotic system. It is noted in [15] that chaotic test problems tend to exaggerate differences between methods. The initial conditions are $(X_0, Y_0, Z_0) = (0, 1, 0)$. The system is solved over the time interval $[0, 5]$ with backward Euler, backward Euler plus filter and BDF2 with constant time step. A reference solution is obtained by self-adaptive RK4-5. We present solutions of the Lorenz system for $\Delta t = 0.01$ (left) and 0.02 (right) in Figure 4.5. The left figure shows that for moderately small time step backward Euler over-damps severely while both BDF2 and backward Euler plus filter are accurate, even for constant time steps. Both have a small positive phase error. The right figure in Figure 4.5 shows that for large time step, each is inaccurate.

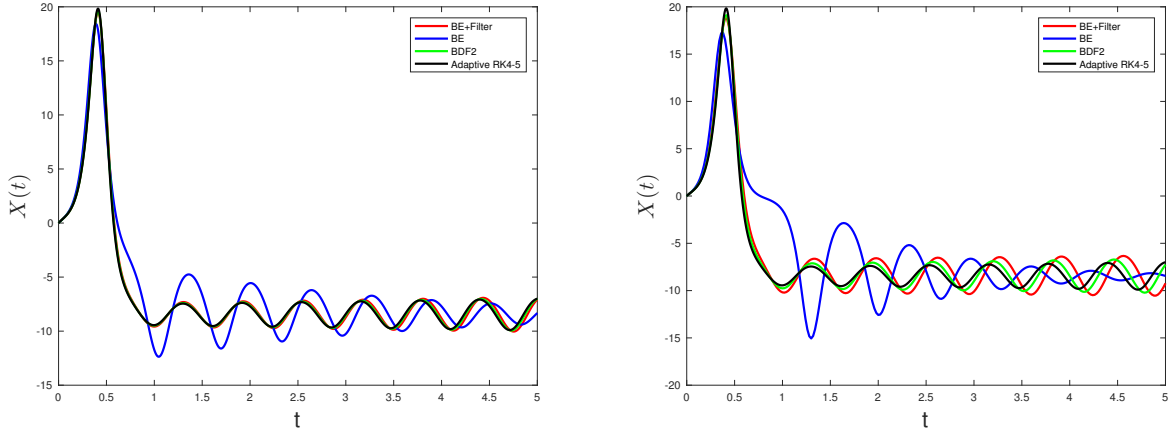


Figure 4.5: Numerical solution of X for the Lorenz system with time step $\Delta t = 0.01$ (left) and $\Delta t = 0.02$ (right).

4.4.2 Preservation of Lyapunov Stability

The implicit method approaches steady state in the test problem below, from [44]. The nonlinear system is

$$\begin{aligned}\frac{du_1}{dt} + u_2 u_2 + u_1 &= 1, \\ \frac{du_2}{dt} - u_2 u_1 + u_2 &= 1.\end{aligned}$$

with $u_1(0) = 0$ and $u_2(0) = 0$. In this test, given in Figure 4.6, adding a filter step preserves Lyapunov stability; the approximate solution (correctly) approaches steady state.

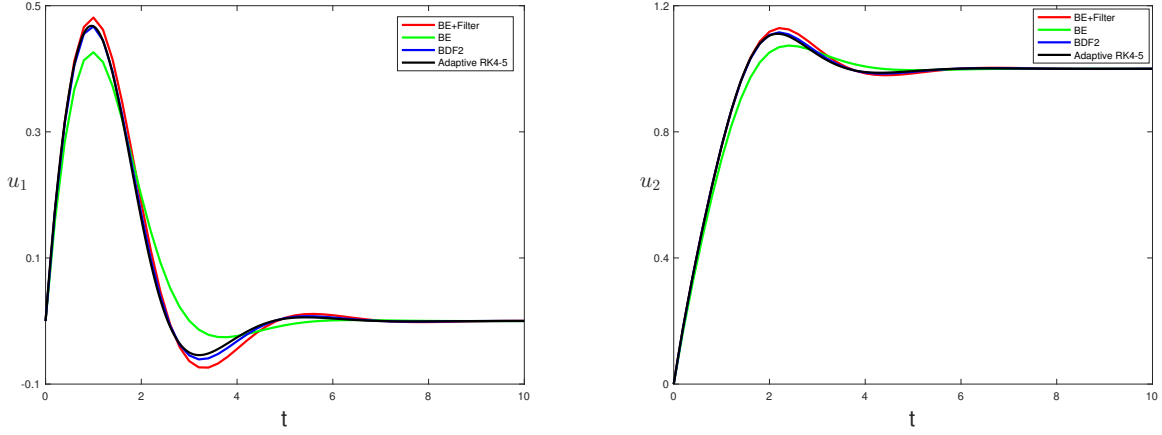


Figure 4.6: Numerical solution of Lyapunov stability test problem with time step $\Delta t = 0.2$

4.4.3 Periodic and Quasi-Periodic Oscillations

Consider the pendulum test problem from [34, 50] given by

$$\begin{aligned}\frac{d\theta}{dt} &= \frac{v}{L}, \\ \frac{dv}{dt} &= -g \sin \theta,\end{aligned}$$

where θ, v, L and g denote, respectively, angular displacement, velocity along the arc, length of the pendulum, and the acceleration due to gravity. Set $\theta(0) = 0.9\pi$, $v(0) = 0$, $g = 9.8$, time step $\Delta t = 0.1$ and $L = 49$. The behavior of the numerical solutions over several periods is depicted in Figure 4.7. Consistently with test 1, the phase and amplitude errors in both backward Euler plus filter and BDF2 are small while both are large for backward Euler. Adding the filter step to backward Euler has greatly increased accuracy.

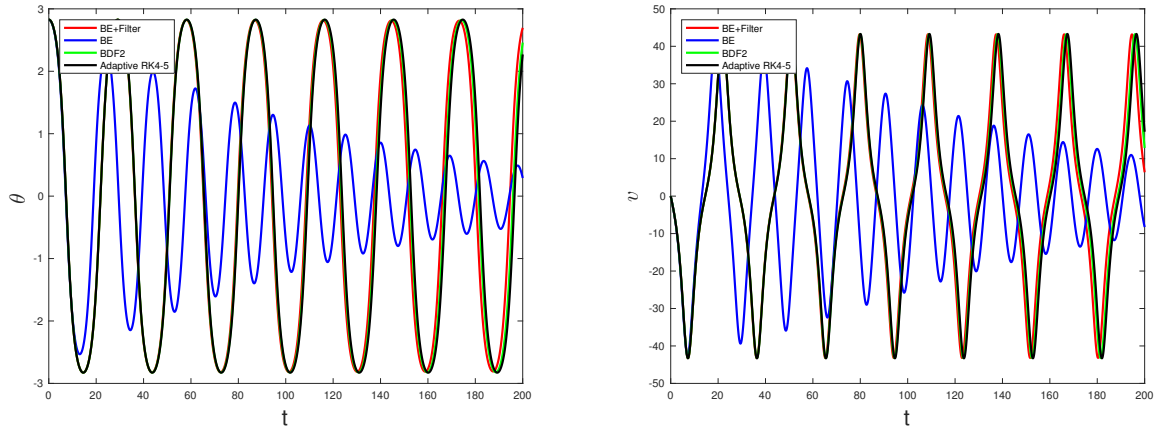


Figure 4.7: Numerical solution of pendulum test problem with time step $\Delta t = 0.1$

Next, we test the method on quasi-periodic oscillations, we solve the following initial value problem(IVP)

$$x'''' + (\pi^2 + 1)x'' + \pi^2 x = 0, \quad 0 < t < 20,$$

$$x(0) = 2, \quad x'(0) = 0, \quad x''(0) = -(1 + \pi^2) \text{ and } x'''(0) = 0,$$

written as a first order system. This has exact solution $x(t) = \cos(t) + \cos(\pi t)$, the sum of two periodic functions with incommensurable periods, hence quasi-periodic, [8]. We solve using backward Euler plus filter with fixed time step $\Delta t = 0.1$ and with a rudimentary adaptive backward Euler plus filter method. In the latter we use initial time step $\Delta t = 0.1$, the heuristic estimator (4.2), tolerance $TOL = 0.1$ and adapt by halving and doubling time step Algorithm 4.1. The plots of both with the exact solution are given in Figure 4.8. This test suggests that quasi-periodic oscillations are a more challenging test than periodic. Adaptivity is required but even simple adaptivity suffices to obtain an accurate solution.

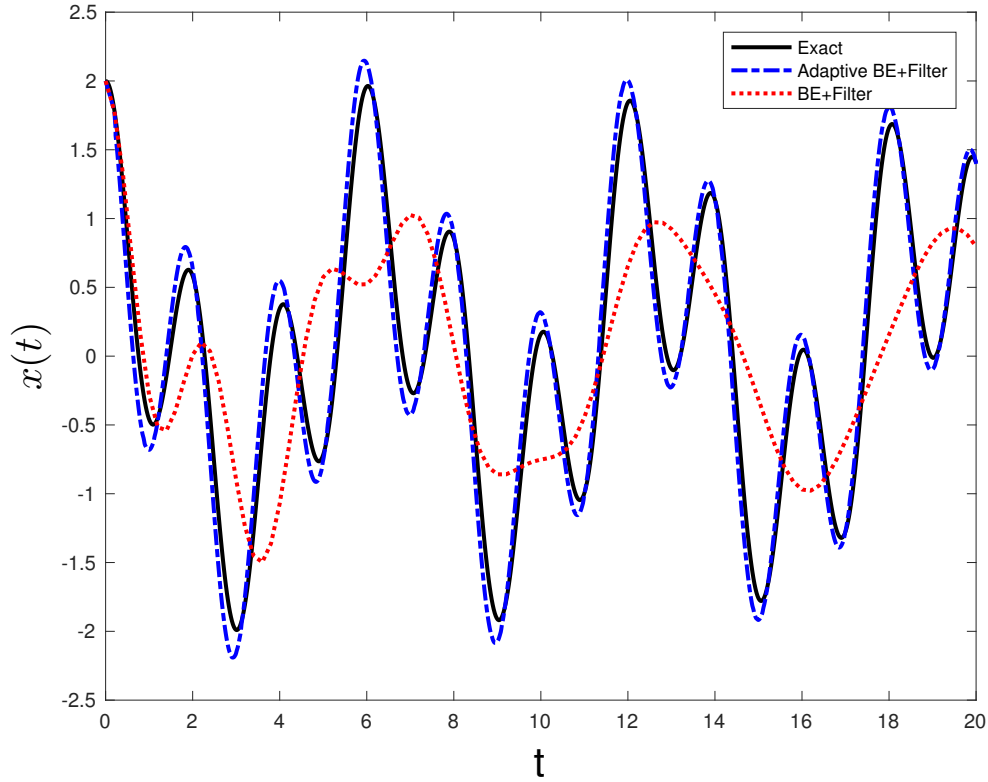


Figure 4.8: Exact soln, non-adaptive and adaptive backward Euler plus filter

4.4.4 The Van der Pol Equation

The last test problem is the Van der Pol equation

$$x'' - \mu(1 - x^2)x' + x = 0,$$

$$x(0) = 2 \text{ and } x'(0) = 0.$$

The Van der Pol equation with parameter $\mu = 1000$ is a common test problem for stiff solvers. The Matlab routine `ode15s` with relative tolerance 10^{-14} and absolute tolerance 10^{-16} provided a reference solution. We solved this with adaptive backward Euler and adaptive backward Euler plus filter for tolerances 10^{-4} and 10^{-6} . The approximate solutions and the time step evolutions are presented in Figure 4.9. The total number of halving, doubling and the same step is given in Table 4.1.

Table 4.1: The comparison of halving, doubling and the same step using variable time step backward Euler and backward Euler plus filter algorithm for the Van der Pol equation.

Method	Halving	Doubling	Same	Tolerance
backward Euler	185	201	41317	10^{-4}
backward Euler plus filter	278	295	7083	10^{-4}
backward Euler	208	228	415519	10^{-6}
backward Euler plus filter	968	990	31830	10^{-6}

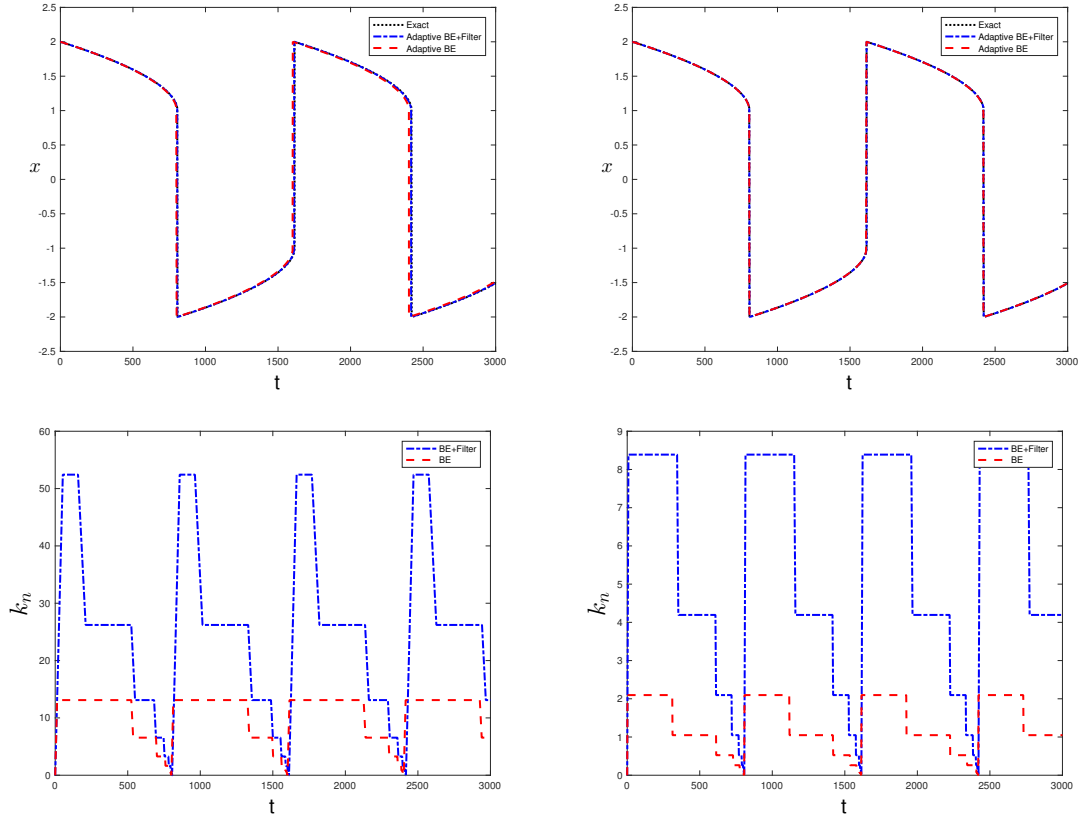


Figure 4.9: Numerical solution of Van der Pol equation using variable time step backward Euler and backward Euler plus filter with 10^{-4} (left) and 10^{-6} (right) tolerances.

4.5 SUMMARY

While a satisfactory, variable time step BDF2 method exists, the combination of backward Euler plus a curvature reducing time filter gives another option that is conceptually clear. Its implementation does require storage of one extra previous value u_{n-1} . The usual barrier in legacy codes is simply identifying the variables needed to be filtered. If this is possible, the filter is easily added by one additional line. If the code solves a partial differential equation, the filter step is a single instruction multiple data (SIMD) instruction across the spacial points. Thus, it maps well to many architectures. Both theory and numerical test show that backward Euler with filter reduces discrete curvature of the solution, increases accuracy from first to second-order, gives an immediate error estimator and induces a method akin to BDF2.

5.0 CONCLUSIONS

The main contribution of this study consists in the construction and analysis of non-intrusive techniques, time filters, which improve the quality of solutions to existing numerical methods, possibly legacy codes.

In Chapter 2 we constructed a higher order Robert-Asselin-Williams filtered leapfrog scheme. The effect of the hoRAW time filter has been analyzed and compared numerically with LF-hoRA and AB3. The hoRAW time filter increases the stability, improves the accuracy of the amplitude of the physical mode up-to two significant digits, effectively suppresses the computational modes, and further diminishes the numerical damping of the hoRA filter. The hoRAW time filter has a twenty percent increase in stability compared hoRA, and the LF-hoRAW is twenty five percent more stable than the AB3 method.

In Chapter 3 we used the notion of modified equations to derive the phase and amplitude errors of pre-defined order for a general high-order Robert-Asselin time filter.

In Chapter 4 we constructed and analyzed a backward Euler plus filter method for constant and variable timesteps. We showed that adding the filter step to the backward Euler method increases accuracy to second-order, reduces oscillations and anti-diffuses the backward Euler method. The numerical tests on oscillatory problems indicate that the numerical dissipation in the new method is comparable to BDF2.

We constructed time filters for two widely used numerical methods for approximation of differential equations, the leapfrog and backward Euler schemes. The combinations of the methods with our time filters are at least second-order accurate and more stable than the original numerical methods, at practically no extra-cost.

BIBLIOGRAPHY

- [1] J. AMEZCUA, E. KALNAY, AND P. D. WILLIAMS, *The effects of the RAW filter on the climatology and forecast skill of the SPEEDY model*, Mon. Wea. Rev., 139 (2011), pp. 608–619.
- [2] R. ASSELIN, *Frequency filter for time integrations*, Mon. Wea. Rev., 100 (1972), pp. 487–490.
- [3] J. BECKER, *A second order backward difference method with variable steps for a parabolic problem*, BIT, 38 (1998), pp. 644–662.
- [4] R. BLECK, *Short-range prediction in isentropic coordinates with filtered and unfiltered numerical models*, Mon. Wea. Rev., 102 (1974), pp. 813–829.
- [5] M. P. CALVO, A. MURUA, AND J. M. SANZ-SERNA, *Modified equations for ODEs*, in Chaotic numerics (Geelong, 1993), vol. 172 of Contemp. Math., Amer. Math. Soc., Providence, RI, 1994, pp. 63–74.
- [6] S. D. CONTE AND C. W. D. BOOR, *Elementary Numerical Analysis: An Algorithmic Approach*, McGraw-Hill Higher Education, 3rd ed., 1980.
- [7] E. CORDERO AND A. STANIFORTH, *A problem with the Robert–Asselin time filter for three-time-level semi-implicit semi-lagrangian discretizations*, Mon. Wea. Rev., 132 (2004), pp. 600–610.
- [8] C. CORDUNEANU AND H. BOHR, *Almost Periodic Functions*, Chelsea, 1989.
- [9] G. DAHLQUIST, *A special stability problem for linear multistep methods*, BIT Numerical Mathematics, 3 (1963), pp. 27–43.
- [10] ———, *Positive functions and some applications to stability questions for numerical methods*, in Recent advances in numerical analysis (Proc. Sympos., Math. Res. Center, Univ. Wisconsin, Madison, Wis., 1978), vol. 41 of Publ. Math. Res. Center Univ. Wisconsin, Academic Press, New York-London, 1978, pp. 1–29.
- [11] ———, *Some contractivity questions for one-leg- and linear multistep methods*, Tech. Rep. TRITA-NA-79-05, Stockholm Univ. Royal Inst. Technol., Stockholm, 1979.

- [12] —, *Some properties of linear multistep and one-leg methods for ordinary differential equations*, Tech. Rep. TRITA-NA-79-04, Stockholm Univ. Royal Inst. Technol., Stockholm, Apr 1979.
- [13] G. G. DAHLQUIST, W. LINIGER, AND O. NEVANLINNA, *Stability of two-step methods for variable integration steps*, SIAM J. Numer. Anal., 20 (1983), pp. 1071–1085.
- [14] R. DALEY, C. GIRARD, J. HENDERSON, AND I. SIMMONDS, *Short-term forecasting with a multilevel spectral primitive equation model part I - model formulation*, Atmosphere, 14 (1976), pp. 98–116.
- [15] D. R. DURRAN, *The third-order Adams-Bashforth method: An attractive alternative to leapfrog time differencing*, Mon. Wea. Rev., 119 (1991), pp. 702–720.
- [16] —, *Numerical methods for fluid dynamics*, vol. 32 of Texts in Applied Mathematics, Springer, New York, second ed., 2010. With applications to geophysics.
- [17] E. EMMRICH, *Stability and error of the variable two-step BDF for semilinear parabolic problems*, J. Appl. Math. Comput., 19 (2005), pp. 33–55.
- [18] C. W. GEAR, *Numerical Initial Value Problems in Ordinary Differential Equations*, Prentice-Hall, Englewood-Cliffs, NJ, 1971.
- [19] D. F. GRIFFITHS AND D. J. HIGHAM, *Numerical methods for ordinary differential equations*, Springer Undergraduate Mathematics Series, Springer-Verlag London, Ltd., London, 2010. Initial value problems.
- [20] D. F. GRIFFITHS AND J. M. SANZ-SERNA, *On the scope of the method of modified equations*, SIAM J. Sci. Statist. Comput., 7 (1986), pp. 994–1008.
- [21] R. D. GRIGORIEFF, *Stability of multistep-methods on variable grids*, Numer. Math., 42 (1983), pp. 359–377.
- [22] A. GUZEL AND W. LAYTON, *Time filters increase accuracy of the fully implicit method*, BIT Numerical Mathematics, (2018).
- [23] A. GUZEL AND C. TRENCH, *The Williams step increases the stability and accuracy of the hoRA time filter*, tech. rep., University of Pittsburgh, 2017.
- [24] E. HAIRER, S. P. NØRSETT, AND G. WANNER, *Solving ordinary differential equations. I*, vol. 8 of Springer Series in Computational Mathematics, Springer-Verlag, Berlin, second ed., 1993. Nonstiff problems.
- [25] E. HAIRER AND G. WANNER, *Solving ordinary differential equations. II*, vol. 14 of Springer Series in Computational Mathematics, Springer-Verlag, Berlin, 2010. Stiff and differential-algebraic problems, Second revised edition.

- [26] N. J. HIGHAM, *Accuracy and stability of numerical algorithms*, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1996.
- [27] W. HUNDSDOERFER AND J. VERWER, *Numerical solution of time-dependent advection-diffusion-reaction equations*, vol. 33 of Springer Series in Computational Mathematics, Springer-Verlag, Berlin, 2003.
- [28] N. HURL, W. LAYTON, Y. LI, AND C. TRENCH, *Stability analysis of the Crank–Nicolson–Leapfrog method with the Robert–Asselin–Williams time filter*, BIT, 54 (2014), pp. 1009–1021.
- [29] E. KALNAY, *Atmospheric Modeling, Data Assimilation and Predictability*, Cambridge University Press, 2003.
- [30] M. KWIZAK AND A. J. ROBERT, *A semi-implicit scheme for grid point atmospheric models of the primitive equation*, Mon. Wea. Rev., 99 (1971), pp. 32–36.
- [31] W. LAYTON, Y. LI, AND C. TRENCH, *Recent developments in IMEX methods with time filters for systems of evolution equations*, J. Comput. Appl. Math., 299 (2016), pp. 50–67.
- [32] Y. LI, *Time filters for numerical weather prediction*, PhD thesis, University of Pittsburgh, 2016.
- [33] Y. LI AND C. TRENCH, *A higher-order Robert–Asselin type time filter*, J. Comput. Phys., 259 (2014), pp. 23–32.
- [34] ———, *Analysis of time filters used with the leapfrog scheme*, in Coupled Problems in Science and Engineering VI COUPLED PROBLEMS 2015, E. O. n. B. Schrefler and M. Papadrakakis, eds., Barcelona, Spain, april 2015, International Center for Numerical Methods in Engineering (CIMNE), pp. 1261–1272.
- [35] E. N. LORENZ, *Deterministic nonperiodic flow*, Journal of the atmospheric sciences, 20 (1963), pp. 130–141.
- [36] L. NAJMAN AND P. ROMON, *Modern Approaches to Discrete Curvature*, vol. 2184 of Lecture Note in Mathematics, Springer International Publishing, 2017.
- [37] O. NEVANLINNA, *Some remarks on variable step integration*, Z. Angew. Math. Mech., 64 (1984), pp. 315–316.
- [38] N. OGER, O. PANNEKOUCKE, A. DOERENBECHER, AND P. ARBOGAST, *Assessing the influence of the model trajectory in the adaptive observation Kalman filter sensitivity method*, Q. J. R. Meteorol. Soc, 138 (2012), pp. 813–825.
- [39] A. QUARTERONI, R. SACCO, AND F. SALERI, *Numerical mathematics*, vol. 37 of Texts in Applied Mathematics, Springer-Verlag, Berlin, second ed., 2007.

- [40] D. REN AND L. M. LESLIE, *Three positive feedback mechanisms for ice-sheet melting in a warming climate*, J. Glaciol., 57 (2011), pp. 1057–1066.
- [41] A. ROBERT AND M. LÉPINE, *An anomaly in the behaviour of the time filter used with the leapfrog scheme in atmospheric models*, Atmosphere-Ocean, 35 (1997), pp. S3–S15.
- [42] A. J. ROBERT, *The integration of a spectral model of the atmosphere by the implicit method*, Proc. WMO-IUGG Symp. on NWP, Tokyo, Japan Meteorological Agency, (1969), pp. 19–24.
- [43] L. F. SHAMPINE, *Error estimation and control for odes*, Journal of Scientific Computing, 25 (2005), pp. 3–16.
- [44] M. SUSSMAN, *A stability example*, tech. rep., University of Pittsburgh, 2010.
- [45] B. VAN DER POL D.SC. AND J. VAN DER MARK, *The heartbeat considered as a relaxation oscillation, and an electrical model of the heart*, The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science, 6 (1928), pp. 763–775.
- [46] R. F. WARMING AND B. J. HYETT, *The modified equation approach to the stability and accuracy analysis of finite-difference methods*, J. Computational Phys., 14 (1974), pp. 159–179.
- [47] M. WATANABE, T. SUZUKI, R. OISHI, Y. KOMURO, S. WATANABE, S. EMORI, T. TAKEMURA, M. CHIKIRA, T. OGURA, M. SEKIGUCHI, K. TAKATA, D. YAMAZAKI, T. YOKOHATA, T. NOZAWA, H. HASUMI, H. TATEBE, AND M. KIMOTO, *Improved climate simulation by MIROC5: Mean states, variability, and climate sensitivity*, Journal of Climate, 23 (2010), pp. 6312–6335.
- [48] P. D. WILLIAMS, *A proposed modification to the Robert–Asselin time filter*, Mon. Wea. Rev., 137 (2009), pp. 2538–2546.
- [49] —, *The RAW filter: An improvement to the Robert–Asselin filter in semi-implicit integrations*, Mon. Wea. Rev., 139 (2011), pp. 1996–2007.
- [50] —, *Achieving seventh-order amplitude accuracy in leapfrog integrations*, Mon. Wea. Rev., 141 (2013), pp. 3037–3051.
- [51] C.-C. YOUNG, Y.-C. LIANG, Y.-H. TSENG, AND C.-H. CHOW, *Characteristics of the RAW-filtered leapfrog time-stepping scheme in the ocean general circulation model*, Mon. Wea. Rev., 142 (2013), pp. 434–447.
- [52] C.-C. YOUNG, Y.-H. TSENG, M.-L. SHEN, Y.-C. LIANG, M.-H. CHEN, AND C.-H. CHIEN, *Software development of the Taiwan Multi-scale Community Ocean Model (TIMCOM)*, Environ. Modell. Software, 38 (2012), pp. 214 – 219.